



A Planning Tool Supporting the Deployment of Cloud Applications

Tudor Alexandru Lascu, Jacopo Mauro, Gianluigi Zavattaro

► **To cite this version:**

Tudor Alexandru Lascu, Jacopo Mauro, Gianluigi Zavattaro. A Planning Tool Supporting the Deployment of Cloud Applications. Other publication. 2013. <hal-00843925>

HAL Id: hal-00843925

<https://hal.inria.fr/hal-00843925>

Submitted on 12 Jul 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Planning Tool Supporting the Deployment of Cloud Applications ¹

Technical Report

Tudor A. Lascu
lascu@cs.unibo.it

Jacopo Mauro
jmauro@cs.unibo.it

Gianluigi Zavattaro
zavattar@cs.unibo.it

June 1st 2013

¹This work was supported by the French ANR project ANR-2010-SEGI-013-01 Aeolus and partially performed at IRILL, center for Free Software Research and Innovation in Paris, France, <http://www.irill.org>

Contents

1 Motivation	2
2 The Deployment Problem	5
3 The Planning Algorithm	7
3.1 Reachability analysis	8
3.2 Abstract Planning	9
3.3 Plan synthesis	13
4 Validation	15
5 Related work	18
6 Conclusions & future work	20

Abstract

Cloud computing offers the possibility to build sophisticated software systems on virtualized infrastructures at a fraction of the cost necessary just a few years ago. Nevertheless, the deployment of such complex systems is a serious issue due to the large number of involved software packages and services, and to their elaborated interdependencies. In this paper we address the challenge of automatizing this complex deployment process. We first formalize it as a planning problem and observe that standard planning tools can effectively solve it only on small and trivial instances. For this reason, we propose an ad hoc planning technique which we validate by means of a prototype implementation able to effectively solve this deployment problem also on instances of realistic size.

Chapter 1

Motivation

Current technology makes it possible to run distributed software systems on-demand on a virtualized infrastructure, at a fraction of the cost which was necessary just a few years ago, thus fulfilling one of the many promises of what is generally referred to as “*cloud computing*”. But realizing, maintaining and reconfiguring such software systems is a serious challenge. When the number of components grows it becomes essential to be able to specify at a certain level of abstraction a particular configuration of the distributed software system, and to rely on tools that provide a set of possible deployment and configuration actions leading to a system configuration that satisfies some user requests. Recent works have introduced formalisms which focus on this *automation* aspect of the deployment process, like the Juju initiative within Ubuntu [19], the Engage system [16], and the Aeolus model [14]. All of them have the goal to provide a planner to enable the automation of finding possible deployment paths. In particular, the Aeolus component model has been designed in order to investigate the computational boundaries of this specific kind of planning problem. Namely, components describe resources which provide and require different functionalities by means of ports, and that may be in conflict with each other. These components are equipped with *state machines* that describe *declaratively* how required and provided functionalities are enacted. Automating a deployment consists of planning a sequence of low-level actions like creation/deletion of components, port binding/unbinding, and internal state changes, that reaches a configuration with at least one component in a specific target internal state. In the full Aeolus model it is also possible to specify *capacity constraints*, i.e. for each provided functionality how many requirements it can satisfy, and for each requirement how many different instances of the provided functionality are needed to satisfy it. In [14] we have proved that the deployment problem is undecidable for the full Aeolus model. On the contrary, if capacity constraints are not considered, we have proved in [13] that the problem turns out to be decidable, but it is Exp-Space hard.

In this work, we focus on the feasibility of developing effective tools able to tackle this problem, and for this reason we have decided to further simplify the Aeolus model by removing also conflicts. Automatic deployment tools like Juju and Engage abstract away from conflicts too. The absence of conflicts is useful, for instance in Engage, to complete partial configurations simply by adding new components without having to check whether these are incompatible with already present components.

In Engage, once the configuration has been completed, a low-level deployment plan

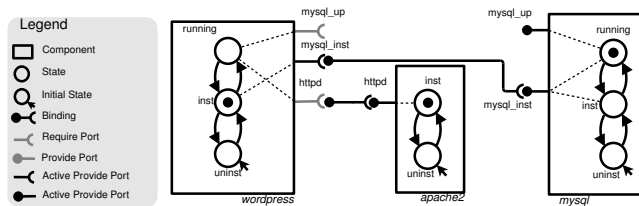


Figure 1.1: Wordpress installation in the Aeolus model.

is synthesized by topologically sorting a graph-based representation of the components in the configuration, where arcs represent dependencies between components. This approach works only if there are no circular dependencies among components. Unfortunately, this form of circularities frequently occur both at the level of fine-grained software packages and of coarse-grained services. For instance, as far as packages are concerned, circular dependencies have to be dealt with for bootstrapping a Linux distribution (see for instance [11] for a list of circular build dependencies between packages in Debian). Concerning services, another example of circular dependency can be found in the master-master replication pattern of `mysql` [9]. In this case, for reliability purposes, two `mysql` servers act at the same time both as master and slave of each other. In the Aeolus model, circular dependencies among components can be naturally represented (see, for instance, the example depicted in Fig. 3.1).

In this work, we present a novel planning technique for the Aeolus model (without conflicts) able to synthesize a deployment plan, when it exists, or to terminate by indicating failure, otherwise. This new technique is based on three distinct phases. In a first phase the existence of a plan is checked by performing a forward symbolic *reachability analysis* of all possible reachable states of the components. If the target state is reachable, a second phase of *abstract planning* generates a graph that indicates the kinds of internal state change actions that are necessary, and the causal dependencies among them. Causal dependencies reflect, for instance, the fact that a component should enter a state enacting a functionality before another component enters a state requiring such functionality. In the third phase of *plan generation* an *adaptive* topological sort of the abstract plan is performed. By adaptive, we mean that the abstract plan could be rearranged during the topological sort if component duplication is needed. Component duplication is the technique used to deal with those cases in which more instances of the same kind of component must be contemporaneously deployed, in different states, in order to enact different functionalities at the same time.

The validation of our planning technique has been achieved by means of a prototype implementation. We compared the prototype on some use cases against some standard general purpose planners providing an encoding of the problem in the Planning Domain Definition Language (PDDL) [17]. Results indicate that general purpose planners can not be used in practice since they do not efficiently handle scenarios involving just 4 or 5 components. Our prototype, instead, produces plans involving hundreds of components in less than one minute.

Structure of the document. In Chapter 2 we recall the Aeolus component model

and we report the corresponding formalization of the component deployment problem as a specific instance of planning. In Chapter 3 we present our algorithm for solving this planning problem while in Chapter 4 we discuss its implementation and validation against standard planning tools. Finally, in Chapter 5 we discuss related works while in Chapter 6 we report some concluding remarks and future directions of development.

Chapter 2

The Deployment Problem

In this chapter we report an informal description of the fragment, without capacity constraints and conflicts, of the Aeolus models formally defined in [14]. Based on the Aeolus representation of the universe of available components, we will be able to formalize the problem of deploying a correct and complete component configuration as a specific case of planning problem.

The Aeolus model has been designed to uniformly represent both fine-grained software components, like packages to be installed on one single virtual machine, and coarse-grained services possibly obtained as composition of distributed and properly connected sub-services. A service, a package or a software component can be modeled as grey-box exposing (only) relevant behaviour by means of internal states and actions for changing them during deployment. For instance Fig. 1.1 depicts three components representing, respectively, *wordpress*, *apache2*, and *mysql* services : *wordpress* and *apache2* are installed while *mysql* is in running mode.

In Aeolus components are modeled as finite state automata indicating the possible internal states and the corresponding state transitions. Each state may activate provide and require ports: active require ports represent the functionalities needed by the service to be in a given state and must be bound to active provide ports of other components. For instance, *mysql* can go from the *running* to an *inst* state. Doing this it will no longer provide the port with interface *mysql_up* but just the port with *mysql_inst* interface. On the other hand *wordpress* can not switch its state to *run* because in order to do so it requires the port *mysql_up* that is not yet bound to an active provide port. Ports are labeled with *interface* names indicating the corresponding required/provided functionality. We can see in the given example that *wordpress* to be installed requires at least an installation of *mysql*. This is expressed by the fact that the *inst* state of *wordpress* requires a port with interface *mysql_inst* that is provided by *mysql* both in *running* and *installed* mode. Initial states, corresponding to state *uninstalled*, do not have any requirement.

According to the Aeolus model, the deployment problem is formalized as follows. The input of the problem is:

- a *universe* of possible component types described according to the approach reported above;
- a *target state*, i.e. an internal state of one specific component type that must be active in the configuration to be deployed.

The output of the problem is a *sequence of low-level deployment actions* on components of the type in the input universe that, upon execution, produces a configuration containing at least one component in the target state. Moreover, all the intermediary configurations, as well as the finally reached one, must be correct in the sense that all active require ports should be bound to complementary active provide ports of other components.

The possible low-level actions are as follows:

- $create(C, id)$ that creates a new component of type C in an initial state. The new component is identified by a unique and fresh identifier id ;
- $delete(id)$ that deletes the component identified by the identifier id ;
- $bind(r, id_1, id_2)$ that creates a binding between the provided port r of the component identified by id_1 and the required port of the component identified by id_2 ;
- $unbind(r, id_1, id_2)$ that deletes the binding between the provided port r of the component identified by id_1 and the required port of the component identified by id_2 ;
- $stateChange(id, s_0, s_1)$ that changes the state of the component identified by id from s_0 to s_1 .

It is worth noticing that there can be more than one way to reach a given configuration of components. For instance, one possible way to obtain the configuration depicted in Fig. 1.1 from scratch, is to first create the resources via the actions $create(wordpress, w)$, $create(apache2, a)$, and $create(mysql, m)$. These three actions create three new components identified by w , a , and m respectively. All these new components will be in the *uninst* state that is the initial state for all of them. Then the *apache2* and *mysql* components can be installed by performing the action $stateChange(a, uninst, inst)$ and $stateChange(m, uninst, inst)$. At this point, to be able to install *wordpress*, we need first to bind the *mysql_inst* port. This is done by performing $bind(mysql_inst, m, w)$. After the creation of the binding, *wordpress* can be installed by performing $stateChange(w, uninst, inst)$. Finally the configuration depicted in Fig. 1.1 can be obtained by performing the $bind(httpd, a, w)$ and $stateChange(m, inst, run)$ actions.

Note that the *unbind*, *delete*, and *stateChange* actions sometimes cannot be performed since their execution would violate the constraint that each active require port must be bound to an active provide one. *bind* and *create* actions, instead, can always be performed as bindings are allowed between ports that are not active and we require that initial states do not activate require ports.

As a final remark, we observe that the decision to use one unique internal *target state* to specify the configuration to be reached is not a limitation. In fact, this target state could activate several require ports indicating an entire set of functionalities that must be present in the final configuration.

Chapter 3

The Planning Algorithm

In this chapter we propose a novel technique to solve the deployment problem, defined in previous section. The technique is divided in three phases, namely, *reachability analysis*, *abstract planning* and *plan generation*,

The first phase computes the states of the components that can be obtained, starting from an empty configuration. If the target state can be reached, an abstract plan is generated describing the needed types of components and a path to reach the target state. Subsequently a concrete plan is obtained by specifically instantiating the component types selected in the abstract plan.

As a running example we model the compilation of package kerberos with ldap support in a Debian system. To build kerberos (*krb5*) the *libldap2-dev* package of *openldap* is needed. This package however depends on *libkrb5-dev* from *krb5*. There is therefore a circular dependency between *krb5* and *openldap*. In Debian the generic way to deal with these circular dependencies is profile builds: every package caters for multiple stages of staged/bootstrap build, so that if necessary a package can have *stage1*, *stage2*, ... before the final, *normal*, build. This allows the depending package to be built allowing the installation of the target package with all the desired functionalities. In the kerberos case, *krb5* is built in the first stage missing out the generation of the *krb5-ldap* package. Then *openldap* can be built directly into its *normal* build satisfying its dependencies. Once *openldap* is built, *krb5* can also be build into its *normal* stage. This process would be modeled in Aeolus as depicted in Fig. 3.1.

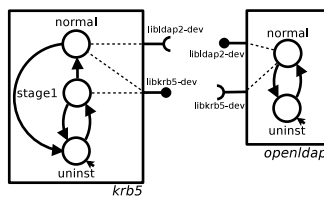


Figure 3.1: Representation of the *krb5* and *openldap* components.

3.1 Reachability analysis

The first step in the proposed technique checks if the the desired target state can be reached.

To do so we rely on an abstraction that ignores bindings and *delete* actions. The abstraction on the bindings is possible since we can safely assume that, given a set of components, every complementary port on two distinct components is bound. Delete actions are not necessary since the presence of one component does not hinder the reachability of one state in another component.¹

All reachable states are computed, for each of the component types in the given universe. In the following we use the pair $\langle \mathcal{T}, q \rangle$ to denote a component type \mathcal{T} and one of its state q .

An increasing sequence of sets of component-state pairs S_0, \dots, S_n is built in such a way that S_{i+1} extends S_i with the new states that can be reached upon execution of a *stateChange* action. The first set, S_0 , contains all the components in their initial state, i.e. $S_0 = \{ \langle \mathcal{T}, q_0 \rangle \mid q_0 \text{ initial state of } \mathcal{T} \}$. Formally S_{i+1} is the largest set satisfying the following constraints:

- $S_i \subseteq S_{i+1}$;
- $\langle \mathcal{T}, q \rangle \in S_{i+1}$ implies the existence of $\langle \mathcal{T}, q' \rangle \in S_i$ such that there is a transition from q' to q in the state automaton of \mathcal{T} ;
- $\langle \mathcal{T}, q \rangle \in S_{i+1}$ implies that for every require port r activated by the state q of \mathcal{T} there exists $\langle \mathcal{T}', q' \rangle \in S_i$ such that the state q' of \mathcal{T}' activates a provide port r .

The generation of sets proceeds until a fix-point is reached (i.e. $S_{i+1} = S_i$) or the target state is obtained. When the fix-point is reached, if the last set does not contain the target pair it means a plan to achieve the goal does not exist and therefore the procedure terminates. Otherwise, we continue with the next phase.

As input to the next phase, we consider a graph-like representation, called *reachability graph*, of the sets S_0, \dots, S_n that keeps track of all the possible ways to obtain the component state-pairs at level $i+1$ from those at level i . More precisely, the graph has as nodes the pairs in S_0, \dots, S_n : if one node at level $i+1$ was already present at level i , the two nodes are connected with an arc $\cdots\cdots$, if a state pair $\langle \mathcal{T}, q \rangle$ at level $i+1$ can be obtained from $\langle \mathcal{T}, q' \rangle$ at level i by means of a *stateChange* action, an \longrightarrow arc from the former to the latter is added. Visually the reachability graph can therefore be seen as a pyramid of levels of component-states having arrows \longrightarrow or arcs $\cdots\cdots$ between two consecutive levels. Figure 3.2, for instance, depicts the pyramid obtained considering the Kerberos use case.

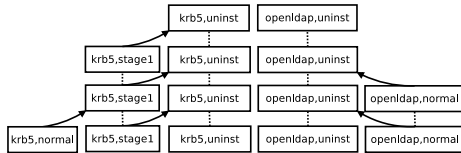


Figure 3.2: Reachability graph for the kerberos running example.

¹This is not true if conflicts are considered.

The first level of Fig. 3.2 contains the two components *krb5* and *openldap* in their initial states. In the second level the component *krb5* in *stage1* state is added since it can be derived from the *krb5* component in state *uninst*. The component *openldap* in *normal* state can not be added at this level since it requires the interface *libkrb5-dev*, not yet provided. *openldap* in *normal* state is added however in the third level since *libkrb5-dev* is now provided by *krb5* in state *stage1*. Finally, in the fourth level, the target state is added deriving it from *krb5* in state *stage1*. This last level is also the fix-point since no new component-state pairs can be generated from it.

Note that keeping component copies allows one to consider different ways a component can use to reach a state. This adds flexibility in deciding how a target can be reached.

In Fig. 1 the pseudo-code for the generation of the reachability graph, is given.

Algorithm 1 Reachability graph construction

```

1: function BUILDREACHABILITYGRAPH( $U, \mathcal{T}_{target}, q_{target}$ )
2:    $Nodes_0 := \{ \langle \mathcal{T}, \mathcal{T}.init \rangle \mid \mathcal{T} \in U \}$ 
3:    $provPort := \bigcup_{\langle \mathcal{T}, q \rangle \in Nodes_0} \{ \mathcal{T}.P(q) \}$ 
4:    $i := 1$ 
5:   repeat
6:      $newNodes := \emptyset$ 
7:      $Arcs_i := \emptyset$ 
8:     for all  $\langle \mathcal{T}, q \rangle \in Nodes_{i-1}$  do
9:       for all  $(q, q') \in \mathcal{T}.trans$  do
10:        if  $\mathcal{T}.R(q') \subseteq provPort$  then
11:           $newNodes.add(\langle \mathcal{T}, q' \rangle)$ 
12:        for all  $\langle \mathcal{T}, q \rangle \in newNodes$  do
13:           $provPort.add(\mathcal{T}.P(q))$ 
14:         $Nodes_i := Nodes_{i-1} \cup newNodes$ 
15:        for all  $\langle \mathcal{T}, q \rangle \in Nodes_{i-1}, \langle \mathcal{T}, q' \rangle \in Nodes_i$  do
16:          if  $(q, q') \in \mathcal{T}.trans \vee (q = q')$  then
17:             $Arcs_i = Arcs_i \cup \{ \langle \langle \mathcal{T}, q' \rangle, \langle \mathcal{T}, q \rangle \rangle \}$ 
18:          if  $(q, q') \in \mathcal{T}.trans$  then
19:             $Arcs_i := Arcs_i \cup \{ \langle \mathcal{T}, q' \rangle \longrightarrow \langle \mathcal{T}, q \rangle \}$ 
20:          if  $q = q'$  then
21:             $Arcs_i := Arcs_i \cup \{ \langle \mathcal{T}, q' \rangle \cdots \langle \mathcal{T}, q \rangle \}$ 
22:           $i := i + 1$ 
23:   until  $\langle \mathcal{T}_{target}, q_{target} \rangle \in Nodes_i \vee newNodes = Nodes_i$ 

```

3.2 Abstract Planning

The reachability graph, previously described, is then used to compute an *abstract plan*. As in the case of the reachability analysis, we abstract from the bind and unbind actions considering instances derived by the reachability analysis pyramid and the temporal dependencies between the creation, deletion and state changes of these instances. We first describe the structure of an abstract plan and then explain how this can be derived from the reachability graph. An abstract plan is a directed graph where the nodes represent

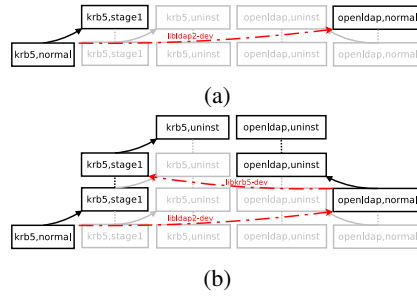


Figure 3.3: Generation of abstract plan for the kerberos example.

either a *create*, *delete*, or *stateChange* action, and arcs represent action precedence constraints. In the following we denote with $\langle z, q, q' \rangle$ a *stateChange* from q to q' of instance z , with $\langle z, \varepsilon, q_0 \rangle$ the *create* action of the instance z in the initial state q_0 , and with $\langle z, q, \varepsilon \rangle$ the *delete* action on the instance z in state q . We consider three types of precedence arcs:

- \longrightarrow : states the precedence of *stateChange* actions on the same component instance; formally $\langle z, x, x' \rangle \longrightarrow \langle z, x', x'' \rangle$ where x' is a state and x, x'' are either states or the special symbol ε denoting absence of the instance z ;
- \xrightarrow{r} where r is an interface: states that if an action deploys an instance z' in a state y' requiring r , provided by z in state y , then state y must be entered before entering state y' , formally $\langle z, x, y \rangle \xrightarrow{r} \langle z', x', y' \rangle$;
- \xrightarrow{r} , where r is an interface, is the dual of the previous arrow: it states that if an action deploys an instance z' in a state y' requiring r , provided by z in state y , then state y' must be exited before exiting state y , formally $\langle z', y', u' \rangle \xrightarrow{r} \langle z, y, u \rangle$.

We are now ready to describe how an abstract plan is obtained. Starting from the reachability graph we select the target component-state pair at the bottom of the pyramid. From the bottom level we then proceed upward selecting the components that are used to deploy the selected component-state pairs at the lower level. To do so, for every selected component at level $i + 1$, we select at level i one of its predecessors (i.e. a component-state pair connected via the \longrightarrow arrow) or a copy (i.e. a component-state pair connected via the \cdots arc). Moreover, for every require port activated by the selected component-state pairs of level $i + 1$ that are not copies, we select a component-state pair at level i that is able to satisfy the requirement, and we keep track of this choice.

In Fig. 3.3a for instance a path from the component `krb5` in state `normal` is selected visiting the component `krb5` in `uninst` state and `inst` state. As can be easily seen a path in the pyramid corresponds to an instance that is created and undergoes some state changes.

To perform these state changes however requirements must be satisfied and therefore new resources may be required to be deployed. Starting from the bottom level we consider every component-state in a path that is not a replica (i.e. obtained using the duplication arcs \cdots) and for every of its required interfaces we select a component-state pair in the previous level that can satisfy its requirements.

For instance in Fig. 3.3b the component `open-ldap` is selected because the component `krb5` in normal state requires the interface `libldap2-dev`.

Then starting from these newly added components we compute a path to reach the top level of the pyramid like it was done for the target component. This process is repeated iteratively until all the non replica component-pair requirements are satisfied by at least another component-pair.

Continuing with our example, in Fig. 3.3b, we select also a path that reaches the component `open-ldap` in normal state. No component need to be selected since the requirements of every component in the path are satisfied (the requirement of `open-ldap` in state normal is satisfied by `krb5` in stage1 state).

We would like to underline that during the selection of component-state pairs different choices could be made. For instance in Fig. 3.3b at the second level we could have selected component `krb5` in state `uninst` to deploy the same component in state `stage1` and component `krb5` in state `stage1` to provide the `libkrb5-dev` interface.

These choices, as we'll see later, have an impact on the number of instances employed to reach the goal. In order to minimize this number we rely on heuristics. In fact, finding the best choice means computing a (global) minimum and would thus imply an exponential cost, due to all alternative possibilities that correspond to disjunctions of conditions.

When we have to select the predecessor of a component, for tie breaking we use as heuristic the following one:

- choose the one that is able to satisfy the maximum number of (not already satisfied) requirements;
- choose the one that can be obtained from an initial configuration satisfying less requirements;
- prefer a copy to one obtained by a state change;
- the one that can be obtained with less state changes from the root.

When we have to select the components needed to satisfy a requirement we use the following strategy:

- choose the one that is able to satisfy the maximum number of (not already satisfied) requirements;
- choose the one that can be obtained from an initial configuration satisfying less requirements;
- the one that can be obtained with less state changes from the root.

Once all the component-state pairs have been selected, we consider a component instance for every maximal path that starts from a component-state in the top level and reaches a component-state that is not a copy. For instance in the kerberos case there are two maximal paths, one starting from the component `krb5` in state `uninst` and reaching the state `normal`, and one starting from the component `openldap` in state `uninst` and reaching the state `normal`. We identify the corresponding instances with z and w respectively.

For every instance we add to the abstract plan its `create`, `delete` and `stateChange` actions. Arrows \longrightarrow are added to connect these actions in chronological order (i.e. first the instance creation, the state changes and then the deletion action). The arrows

\xrightarrow{r} and \dashrightarrow^r are instead added between actions of instances requiring and providing an interface r .

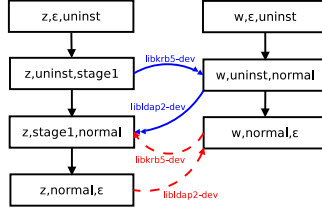


Figure 3.4: Abstract plan for the kerberos running example.

Fig. 3.4 shows the abstract plan obtained for the kerberos case. The four actions on the left are related to instance z while the three on the right are actions related to instance w . z is first created, then it changes its state first into *stage1* and then to *normal* before being deleted. w instead is created, it changes state into *normal*, before being deleted. These precedences are encoded by \rightarrow arrows. z 's requirement of *libldap2-dev* in state *normal*, satisfied by w in *normal* state, is encoded with $\xrightarrow{\text{libldap2-dev}}$ between $\langle w, \text{uninst}, \text{normal} \rangle$ and $\langle z, \text{stage1}, \text{normal} \rangle$ and $\dashrightarrow^{\text{libldap2-dev}}$ between $\langle z, \text{normal}, \epsilon \rangle$ and $\langle w, \text{normal}, \epsilon \rangle$. The first one states that w must be in *normal* state before z moves to *normal* while the second states that the deletion of z must precede the deletion of w . Indeed, if one of these constraints does not hold it means that the abstract plan violates a requirement thus leading to a non correct configuration. Similarly, the *libkrb5-dev* interface requirement of w in state *normal*, satisfied by z in *stage1*, is encoded with $\xrightarrow{\text{libkrb5-dev}}$ between $\langle z, \text{uninst}, \text{stage1} \rangle$ and $\langle w, \text{uninst}, \text{normal} \rangle$ and $\dashrightarrow^{\text{libkrb5-dev}}$ between $\langle w, \text{normal}, \epsilon \rangle$ and $\langle z, \text{stage1}, \text{normal} \rangle$. In this case we can however notice that z continues to provide the port *libkrb5-dev* also when it is in state *normal*. Thus w does not need to be deleted before *krb5* moves to *normal* but it can stay until the *krb5* is not deleted. This relaxation corresponds to setting $\langle z, \text{normal}, \epsilon \rangle$ as the target of $\dashrightarrow^{\text{libkrb5-dev}}$. In general all the constraints $\langle z, x, y \rangle \dashrightarrow^r \langle z, x', y' \rangle$ can be relaxed replacing $\langle z, x', y' \rangle$ with $\langle z, x'', y'' \rangle$ where $\langle z, x'', y'' \rangle$ is a *delete* action or it is the the first *stateChange* reaching a state y'' that does not provide r . After applying these relaxations we obtain the final version of abstract plan that, for the kerberos case, is the one depicted in Fig. 3.5.

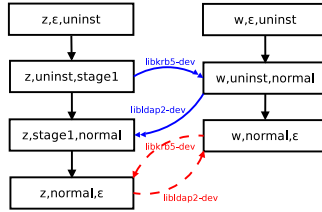


Figure 3.5: Abstract plan for the kerberos example after relaxation.

3.3 Plan synthesis

The abstract plan, computed in previous phase, is then used to synthesize a concrete one. The idea is to visit the nodes of the abstract plan in topological order until the target component is obtained. Visiting a node consists of performing that action. Moreover, in order to properly satisfy component requirements, when an incoming \xrightarrow{r} is encountered a new binding should be created, and when an outgoing \xrightarrow{r} is encountered the corresponding binding should be deleted. Notice that it is not necessary to visit the entire abstract plan as it is sufficient to reach the target state. For this reason, we give priority to the visit of the actions of the components containing such state.

For instance, in the kerberos example, we can extract a concrete plan from the abstract plan in Fig. 3.5 as follows. Assume that the target state is state *normal* of component type *krb5*. As we give priority to the corresponding instance, the first action in the concrete plan is *create(krb5,z)* corresponding to the visit of $\langle z, \varepsilon, uninst \rangle$. The subsequent action is *stateChange(z,uninst,stage1)* corresponding to $\langle z, uninst, stage1 \rangle$. The visit of the actions on the instance *z* cannot proceed due to the incoming arrow $\xrightarrow{libldap2-dev}$; for this reason the next action in the concrete plan is *create(openldap,w)* corresponding to the visit of $\langle w, \varepsilon, uninst \rangle$. The next node in the abstract plan to be visited is $\langle w, uninst, normal \rangle$, but as this node has an incoming $\xrightarrow{libkrb5-dev}$, two actions must be added to the concrete plan: *bind(libkrb5-dev,z,w)* and *stateChange(w,uninst,normal)*. At this point, the visit of the component instance *z* can continue by considering node $\langle z, stage1, normal \rangle$; as this node has an incoming $\xrightarrow{libldap2-dev}$, two actions must be added to the concrete plan: *bind(libldap2-dev,w,z)* and *stateChange(w,uninst,normal)*. This completes the generation of the concrete plan as the target state has been reached.

Unfortunately, the topological visit is not always possible as it may be inhibited by the presence of cycles in the abstract plan. Consider, for instance, a slightly modified version of the kerberos example in which the component type *krb5* in *normal* state requires not only one *openldap* in *normal* state, but also one in *uninst* state. In this case the abstract plan will be as in Fig. 3.6 (note the addition of the pair of arcs labeled with *uninst*).

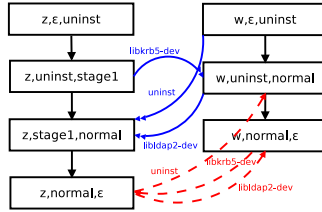


Figure 3.6: Abstract plan for the modified kerberos example.

In the abstract plan in Fig. 3.6 there is a cycle among $\langle w, uninst, normal \rangle$, $\langle w, normal, \varepsilon \rangle$ and $\langle z, normal, \varepsilon \rangle$ that forbids the visit of $\langle w, uninst, normal \rangle$, which is a necessary step to visit the target node $\langle z, stage1, normal \rangle$.

In general, these cycles appear when an instance is expected to provide an interface during a specific phase of the plan, but this is not possible because in the same phase of the plan the instance is required to change its internal state. This problem can be solved by means of *instance duplication*: an additional component instance is

deployed in such a way that the new instance can continue providing the required interface during that specific phase of the plan. The application of instance duplication to the abstract plan in Fig. 3.6 is reported in Fig. 3.7, where we add a new instance y of type *openldap* that does not proceed further than state *uninst*. This new resource is used to satisfy the requirement *uninst* of z . Notice that the topological visit until the target node $\langle z, stage1, normal \rangle$ becomes now possible.

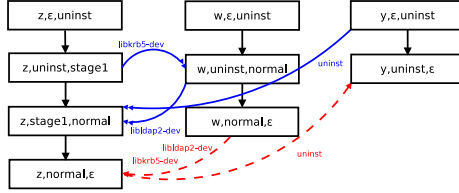


Figure 3.7: Abstract plan in Fig. 3.6 after duplication of instance w .

In general, we have to consider an *adaptive* topological visit, where by adaptive we mean that the abstract plan is transformed when the visit is blocked by a cycle. When a node cannot be visited, we proceed as follows. There is at least an $\overset{r}{\dashrightarrow}$ or a $\overset{r}{\rightarrow}$ arc, incoming to this node. We consider two cases: when there is no incoming $\overset{r}{\rightarrow}$ arc, and otherwise. In the first case, we proceed by removing the incoming $\overset{r}{\dashrightarrow}$ arc, and otherwise. In the first case, we proceed by removing the incoming $\overset{r}{\dashrightarrow}$ by duplicating the instance containing the node that could not be visited. The duplication of the instance is done from its initial node to such node. Concerning the arcs, all the incoming $\overset{r}{\rightarrow}$ and outgoing $\overset{r}{\dashrightarrow}$ arcs of the instance (representing its requirements on the other components) will have to be duplicated and, more important, the incoming $\overset{r}{\dashrightarrow}$ to be removed are moved to the new instance. This is obtained by moving the pairs of outgoing $\overset{r}{\rightarrow}$ / incoming $\overset{r}{\dashrightarrow}$ representing the interface that the new instance must provide to the other components. See, for instance, the pair of arcs labeled with *uninst* that are moved from their initial position in Fig. 3.6 to the new instance y in Fig. 3.7. In this way, the considered node does not any longer have incoming arcs $\overset{r}{\dashrightarrow}$ and it can be visited. Consider now the case in which there are incoming arcs $\overset{r}{\rightarrow}$. We consider the node sources of these arcs. If those nodes cannot be visited, we apply to them the instance duplication procedure. Upon this duplication phase, the initial node could be visited or will have impediments due to incoming arcs $\overset{r}{\dashrightarrow}$. In this case the problem can be solved as described in the previous case.

It is interesting to note that this adaptive visit will eventually terminate because new instances do not introduce new cycles, and because there exists no cycle involving only $\overset{r}{\rightarrow}$ arcs. This is guaranteed by the fact that this class of arcs represents the dependencies selected during the transformation from the reachability graph to the abstract plan; by construction, these dependencies cannot be circular because they always go from nodes at a lower to nodes at a higher level in the pyramid.

Chapter 4

Validation

In the context of knowledge representation and reasoning, a very important application of artificial intelligence, is that of developing languages and tools for reasoning about actions and change and, more specifically, for the problem of planning [8]. Since 1998, a declarative language for planning has been defined for establishing a common syntax for different tools in order to allow different research groups to test their solvers. This language is known as PDDL and his last release is 3.1 (see [?, 17, ?] [17] for information on planning competitions and PDDL).

Our tool solves a planning problem and therefore we tried to validate our ad-hoc planner against standard planners. To do so we have defined an encoding of our specific planning problem into PDDL: each component instance is translated into one PDDL object with possible actions corresponding to state changes. These actions can be acted on the object only when the other objects in the configuration provide the required interfaces. The encoding abstracts from the bind and unbind actions¹ and limits the number of objects that could be concurrently used.²

As a benchmark we have considered Aeolus instances automatically generated following the pattern of component interdependency discussed in Chapter 3 with the kerberos running example. We have considered configurations composed by an increasing number of components having 3 states (excluding one component having only 2 states used to trigger the plan).

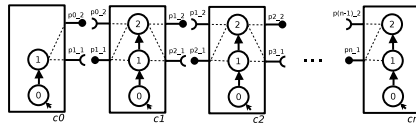


Figure 4.1: Component types considered in our tests.

As depicted in Fig. 4.1, the components have provide and require ports in such a

¹Bind and unbind action can be added to form a valid deployment run in polynomial time in a post processing phase.

²This limitation was necessary because all the solvers assume a finite number of objects –without this limitation the planning problem is undecidable.

Table 4.1 Performances of standard planners.

Size	Test $T+$		Test T	
	Madagascar-p	Metric-FF	Madagascar-p	Metric-FF
3	0.13 s	timeout	0.05 s	0.05 s
4	0.77 s		0.41 s	timeout
5	3.44 s		1.96 s	
6	error		7.22 s	
7			error	

way that a valid plan to reach the state s_2 of the rightmost component would require first to create the components, then perform from cn to c_0 the state change from s_0 to s_1 , and then perform the state change from s_1 to s_2 from c_1 to cn . In these scenarios instance duplication is not needed during the generation of the concrete plan. In order to test also this feature we modified the component types by removing for nearly one fifth of the components the activation of the provide port pX_I from the states s_2 , where X is the number of the selected component. Removing this activation of the provide port requires the duplication of the instance of the component type cX in order to contemporaneously satisfy the requirements of its two neighbor components. The activation of the provide port to be deleted were selected randomly. In the following we will denote respectively with T the test of the scenario where duplication is not needed and with $T+$ the test where duplication is necessary.

The tests were performed using a dual core machine with 2.60 GHz Intel i5 processors, 8GB of RAM, and Ubuntu 13.04 operating system. We used a time cup of 120 seconds and two planners that support the ADL fragment of PDDL (other popular solvers support only fragments of PDDL): Metric-FF [5] and Madagascar-p [3]. The first solver is based on GraphPlan, a standard planning algorithm to prune the search space, while the second encodes the planning problem into a SAT formula and then uses state of the art SAT solvers to solve it. For reducing the search space of these solvers we set to the minimum the number of components that could be used concurrently.

We notice that the performances of the general purpose solvers are quite poor. Table 4.1 reports the time performances for both the T and $T+$ scenarios. With *error* we mean that the solver exited without computing the plan, with *timeout* that the solver took more than 120s and was thus interrupted.

Even though the best solver among the two is Madagascar-p, it is able to solve scenarios with just 6 components for the simpler T test. Metric-FF, in the $T+$ case, times out even when considering a scenario with 3 components. These poor performances are due to the fact that the size of the encoding of the planning problem increases exponentially w.r.t. the number of components that need to be deployed concurrently to reach the target state size. Metric-FF times out because it wastes all its time trying to ground all the possible actions while Madagascar-p fails because the encoding into SAT becomes too big to be handled.

Fig. 4.2 shows instead the time taken by our tool to produce a valid plan for tests T and $T+$. Clearly our ad-hoc planner outperforms the general purpose planners. In

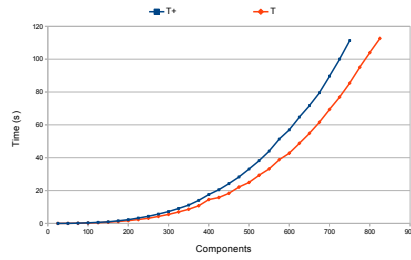


Figure 4.2: Time performances using our tool

both cases it is able to produce plans in less than two minutes for scenarios involving more than 700 components. As expected, duplication affects the performance of the planner, as can be read from comparing the graph relative to $T+$ with the one for T . The degradation of performances, however, is still manageable: plans can be provided in less than a minute even for scenarios with hundreds of components.

The tool was developed in OCaml and is about 3.5 K lines of code. The source code and the problem encoding into the PDDL language are available at [2].

Chapter 5

Related work

The problem of managing networks of interconnected machines has attracted significant attention in the area of system administration. Many popular tools like CFEngine [10], Puppet [20], MCollective [22] and Chef [21] have been developed. Despite their differences, they all allow to declare the components to be installed on each machine, together with their configuration files, and then employ various mechanisms to deploy components accordingly. However, The burden of specifying which components to deploy where, and how to interconnect them is left to the user. These system management tools can however be used as deployment back-end for our planner: once the plan has been computed, the actual deployment can be delegated to this kind of tools. Another system specifically targeted at application deployment in the cloud is Cloud-Foundry [23], but it has the same limitations described above.

The problem of automatizing the selection/distribution/interconnection of the components has been addressed by other tools like ConfSolve [18] and Zephyrus [12]. ConfSolve [18] relies on a constraint solver to propose an optimal allocation of virtual machines on servers and applications on virtual machines, but it does not handle the problem of connections among services. Zephyrus [12] automates assembling complex distributed systems, allowing a user to compute a valid configuration satisfying a high level specification. The description of the model is given in Aeolus, and the tool takes into account conflicts and capacity constraints, but it does not consider components as automata since the only concern is to provide a valid final configuration and not the steps to reach it. Juju [19] and Engage [16], are more similar to our approach: they both rely on a solver to perform the deployment part and they avoid the problem of dealing with conflicts among components. In our approach however we are able to treat circular dependencies that cannot be defined in Engage and should instead be manually resolved using low level programming techniques in Juju. Another approach to automating deployment is proposed in [15]; it uses an Architecture Description Language with information on the relationships among software services, which needs to be explicitly provided by the user in full detail, and uses a decentralized protocol to perform automatic configuration.

Closely related to our work is [7] that proposes an heuristic-based algorithm to remove build dependency cycles for bootstrapping a Linux software distribution. The building order of the packages is generated using a topological sort of a graph. The problem addressed there is, however, slightly but significantly different as the aim is

to build packages from scratch. At every iteration the newer version of the recompiled package replaces the old one. In our case this is not always true as the same service may be required in both installed and running state to support fault tolerance via replication.

Chapter 6

Conclusions & future work

In this work a technique is presented that, given a representation of software packages/services following the Aeolus model, computes the necessary steps to deploy a final configuration. Our solution handles circular dependencies among components and produces plans deploying hundreds of different components in less than a minute. Moreover, the deployment actions can be easily parallelized; once the concrete plan is generated, the knowledge of the precedence constraints among different actions allows one to identify the minimal number of synchronizations needed to perform the deploying actions in parallel.

In order to practically experiment and trial our tool in a real-life environment, with the collaboration of our industrial partner within the ANR project “Aeolus: Mastering Cloud Complexity” [1] –Mandriva SA [4]– we are currently developing a n-tiers application deployment engine. This engine deploys, on an IaaS, a final configuration found by the Zephyrus tool. First, the required virtual machines are provisioned by relying on OpenStack. We then proceed with the installation and configuration phase in which the engine connects to each virtual machine involved in the n-tiers application. The required packages are installed and services are configured following the steps specified by the plan produced by our tool. This is done using the deployment tool MSS (Mandriva Server Setup) [6]. It is worth noticing that these steps are performed without any human interaction, by combining the MSS configuration information, the final configuration computed by Zephyrus and the plan output by our tool. At present we are able to successfully deploy a configuration of several instances of Wordpress backed by a Varnish load balancer. We are currently working on integrating examples with complex database configurations containing master and slave requirements.

In future work, we aim to improve our tool by considering also reconfiguration plans, dealing with cases in which the initial configuration is not empty, but already deployed components could be used and rearranged to reach the desired new configuration. We would also like to take into account conflicts producing a plan that does not violate them or, if this is not feasible in practice, devise plans that minimize the time windows where a system is inconsistent.

Bibliography

- [1] Aeolus: Mastering the Complexity of the Cloud. <http://www.aeolus-project.org/>.
- [2] DACLAMP - Deployment Applications in the CCloud Aeolus Min Planner. <https://github.com/talascu/daclamp>.
- [3] Madagascar-p. <http://users.cecs.anu.edu.au/~jussi/satplan.html>. Retrieved June 2013.
- [4] Mandriva SA. www.mandriva.com.
- [5] Metric-FF. <http://fai.cs.uni-saarland.de/hoffmann/metric-ff.html>. Retrieved June 2013.
- [6] MSS - Mandriva Server Setup. <http://doc.mandriva.com/en/mes5/Enterprise-Server-Manual-EN.html/CS-mss-intro.html>. Retrieved June 2013.
- [7] Pietro Abate and Schauer Johannes. Bootstrapping Software Distributions. In *CBSE*, 2013.
- [8] C. Baral. *Knowledge Representation, Reasoning and Declarative Problem Solving*. Cambridge University Press, 2003.
- [9] Schwartz Baron, Zaitsev Peter, Tkachenko Vadim, Zawodny Jeremy, Lentz Arjen, and Balling Derek J. *High performance mysql, 2nd edition*. O'Reilly, second edition, 2008.
- [10] Mark Burgess. A Site Configuration Engine. *Computing Systems*, 8(2):309–337, 1995.
- [11] Circular Build Dependencies. <http://wiki.debian.org/CircularBuildDependencies>. Retrieved June 2013.
- [12] Roberto Di Cosmo, Michael Lienhardt, Ralf Treinen, Stefano Zacchiroli, and Jakub Zwolakowski. Optimal Provisioning in the Cloud. Technical report, Aeolus project, Juin 2013. <http://hal.archives-ouvertes.fr/hal-00831455>.
- [13] Roberto Di Cosmo, Jacopo Mauro, Stefano Zacchiroli, and Gianluigi Zavattaro. Component Reconfiguration in the Presence of Conflicts. In *ICALP 2013: 40th International Colloquium on Automata, Languages and Programming*, volume 7966 of *LNCS*, pages 187–198. Springer-Verlag, 2013.

- [14] Roberto Di Cosmo, Stefano Zacchiroli, and Gianluigi Zavattaro. Towards a Formal Component Model for the Cloud. In *SEFM 2012: 10th International Conference on Software Engineering and Formal Methods*, volume 7504 of *LNCS*, pages 156–171. Springer-Verlag, 2012.
- [15] X. Etchevers, T. Coupaye, F. Boyer, and N. de Palma. Self-Configuration of Distributed Applications in the Cloud. In *CLOUD 2011*, pages 668–675, july 2011.
- [16] Jeffrey Fischer, Rupak Majumdar, and Shahram Esmaeilsabzali. Engage: A Deployment Management System. In *PLDI'12: Programming Language Design and Implementation*, pages 263–274. ACM, 2012.
- [17] Maria Fox and Derek Long. PDDL2.1: An Extension to PDDL for Expressing Temporal Planning Domains. *J. Artif. Intell. Res. (JAIR)*, 20:61–124, 2003.
- [18] John A. Hewson, Paul Anderson, and Andrew D. Gordon. A Declarative Approach to Automated Configuration. In *LISA '12: Large Installation System Administration Conference*, pages 51–66, 2012.
- [19] Juju, devops distilled. <https://juju.ubuntu.com/>. Retrieved June 2013.
- [20] Luke Kanies. Puppet: Next-generation configuration management. *;login: the USENIX magazine*, 31(1):19–25, 2006.
- [21] Opscode. Chef. <http://www.opscode.com/chef/>. Retrieved June 2013.
- [22] Puppet Labs. Marionette Collective. <http://docs.puppetlabs.com/mcollective/>. Retrieved June 2013.
- [23] VMWare. Cloud Foundry, deploy & scale your applications in seconds. <http://www.cloudfoundry.com/>. Retrieved June 2013.