# Execution and Planning of Distributed Systems Reconfigurations

*ICE Workshop @DisCoTec, 18th of June 2021*
**Hélène Coullon** (Maverick Chardet, Simon Robillard, Dimitri Pertin, Christian Perez, Claude Jard, Didier Lime, Charlène Servantie)

Associate professor at IMT Atlantique, France
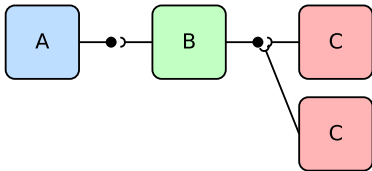
Inria researcher, France

Adjunct professor at UiT, Tromsø, Norway

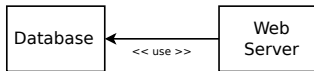# Introduction and motivations

**General definition**

- Non monolithic code,

- modular units of code - **components**,

- software system $=$ **architectural assembly** of component instances,

- interactions between components through **communications**.



- Master/workers,

- microservices,

- service-oriented,

- layered,

- etc.

# Deployment example



Machine 1 **[WHERE]**

**Database (DB) [WHAT]**

**[HOW]** **[LIFECYCLE]**

1. Install
2. Configure
3. Start the service
4. Prepare the service

Machine 2 **[WHERE]**

**Web-server (WS) [WHAT]**

**[HOW]** **[LIFECYCLE]**

1. Install
2. Configure firewall
3. Download
4. Configure parameters
5. Start the service

**[WHEN]****[DEPENDENCIES]**: DB ≪ WS (component granularity)

**[DEPENDENCIES]**: DB(3) ≪ WS(4), DB(4) ≪ WS(5) (lifecycle granularity)

Ever-running and long-running distributed systems
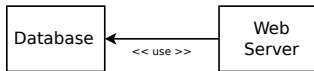
**What is a reconfiguration?**

- Reconfiguration through time
  - need to add/remove components and/or connections
  - need to change internal configurations
- A set of instructions to move from one state of the system to another.

**Examples of reconfiguration reasons**

- Faults or errors on services or hardware (e.g., re-deploy),
- dynamic energy or security constraints (e.g., change the set of components),
- dynamic improvement of performance (e.g., scaling),
- dynamic upgrade of some modules.

**Database (DB) [WHAT]**

**[HOW] [LIFECYCLE]**

1. Backup data
2. Stop the service
3. Download update
4. Configure parameters
5. Start the service
6. Restore data

**Web-server (WS) [WHAT]**

**[HOW] [LIFECYCLE]**

1. Pause the service
2. Configure parameters
3. Start the service

**[DEPENDENCIES]**: WS(1) $\ll$ DB(2), DB(5) $\ll$ WS(2), DB(6) $\ll$ WS(3) (lifecycle granularity)

## Our goals

1. **Efficiency of reconfigurations**
   - reach quickly a targeted configuration,
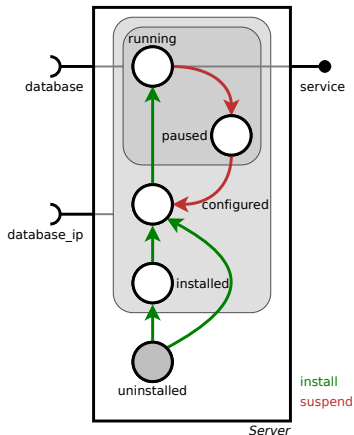   - reduce disruption time.

2. **Execution time prediction**

3. **Safety of reconfigurations**

## Table of Contents

# The Concerto reconfiguration model

# Control components



*Server*

**Written by the component developers**

## Internal net [LIFECYCLE]

- places = milestones
- transitions = concrete actions to perform

## Interfaces [DEPENDENCIES]

- data or service ports
    - use ports = requirements
    - provide ports = provisions
- behaviors
    - subset of transitions
- during execution: active/inactive

## Reconfiguration language

1. Create assemblies of components (software system)

2. Make this assembly evolve at runtime

3. Interact with the lifecycle of components

### Add/remove
Add/remove a component instance to the current assembly

### Connect/disconnect
Connect/disconnect two component instances with compatible ports

### Push behavior
Push a behavior to the behavior queue on a component instance

### Wait
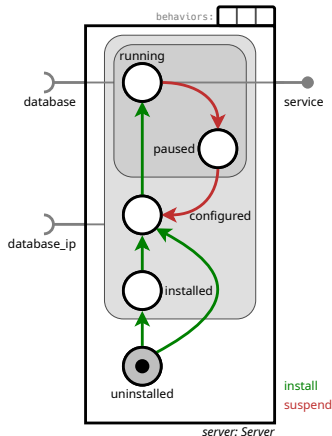Wait for a given component instance or wait all components

Written by the reconfiguration developer

Deployment program:

```
1  add ( server :  Server )
2  add ( db :  Database )
3  con ( server . database_ip , db . ip )
4  con ( server . database , db . service )
5  pushB ( server ,  install )
6  pushB ( db ,  deploy )
7  wait ( server )
```

Written by the reconfiguration developer

Deployment program:

```
1 add(server: Server)
2 add(db: Database)
3 con(server.database_ip,db.ip)
4 con(server.database,db.service)
5 pushB(server, install)
6 pushB(db, deploy)
7 wait(server)
```

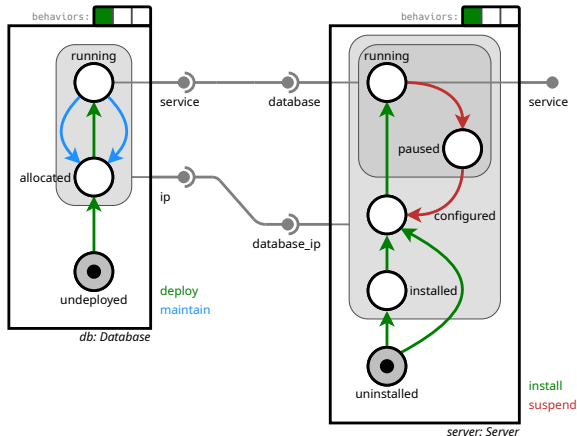Written by the reconfiguration developer

Deployment program:

```
1  add(server: Server)
2  add(db: Database)
3  con(server.database_ip,db.ip)
4  con(server.database,db.service)
5  pushB(server, install)
6  pushB(db, deploy)
7  wait(server)
```
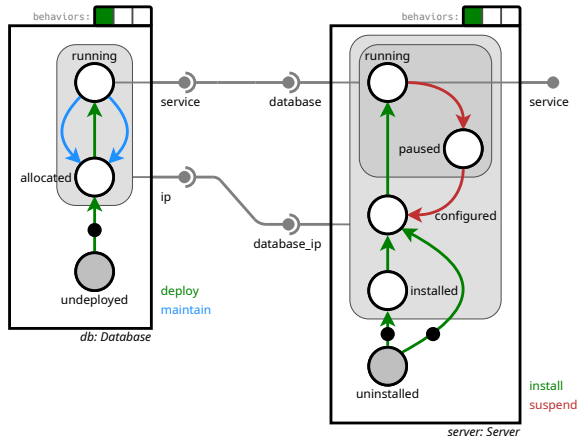
Written by the reconfiguration developer

Deployment program:

```
1 add(server: Server)
2 add(db: Database)
3 con(server.database_ip,db.ip)
4 con(server.database,db.service)
5 pushB(server, install)
6 pushB(db, deploy)
7 wait(server)
```
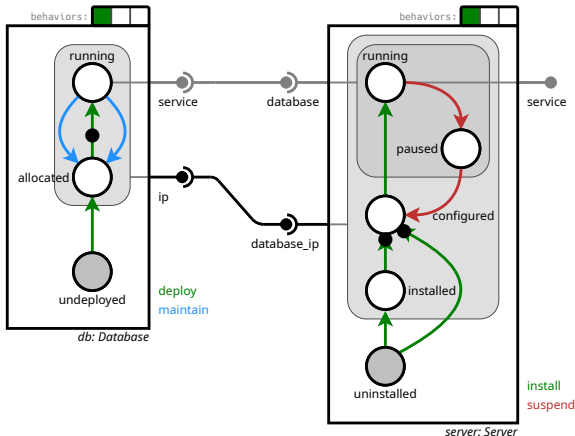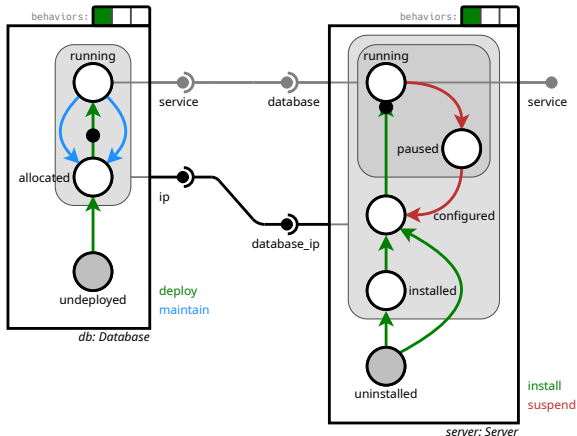
Written by the reconfiguration developer
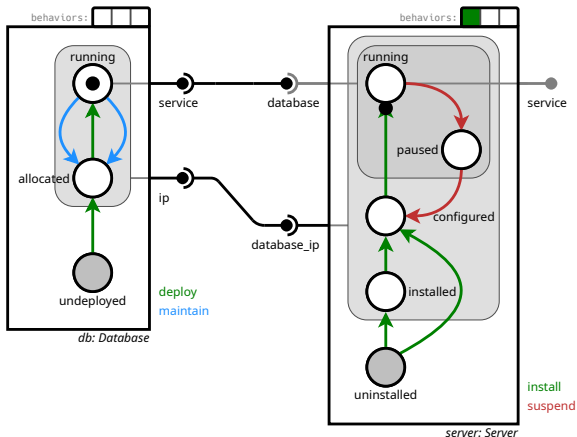
Deployment program:

```
1 add ( server : Server )
2 add ( db : Database )
3 con ( server . database_ip , db . ip )
4 con ( server . database , db . service )
5 pushB ( server , install )
6 pushB ( db , deploy )
7 wait ( server )
```

Written by the reconfiguration developer

Deployment program:

```
1  add ( server :  Server )
2  add ( db :  Database )
3  con ( server . database_ip , db . ip )
4  con ( server . database , db . service )
5  pushB ( server ,  install )
6  pushB ( db ,  deploy )
7  wait ( server )
```

Written by the reconfiguration developer

Deployment program:

```
1 add ( server : Server )
2 add ( db : Database )
3 con ( server . database_ip , db . ip )
4 con ( server . database , db . service )
5 pushB ( server , install )
6 pushB ( db , deploy )
7 wait ( server )
```

Written by the reconfiguration developer

Maintenance program:

```
1  pushB ( db , maintain )
2  pushB ( db , deploy )
3  pushB ( server , suspend )
4  pushB ( server , install )
5  wait ( server )
```

Written by the reconfiguration developer

Maintenance program:

```
1 pushB ( db , maintain )
2 pushB ( db , deploy )
3 pushB ( server , suspend )
4 pushB ( server , install )
5 wait ( server )
```

Written by the reconfiguration developer

Maintenance program:

```
1 pushB(db, maintain)
2 pushB(db, deploy)
3 pushB(server, suspend)
4 pushB(server, install)
5 wait(server)
```
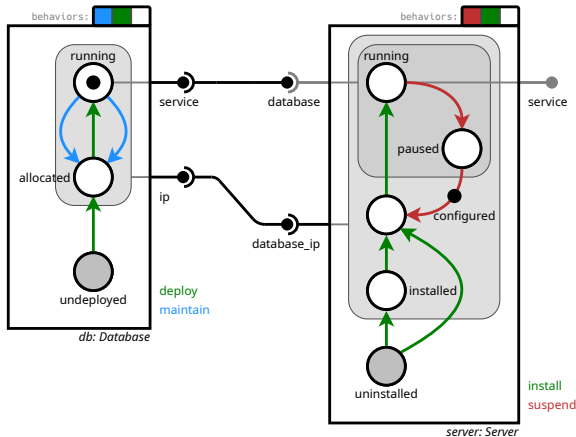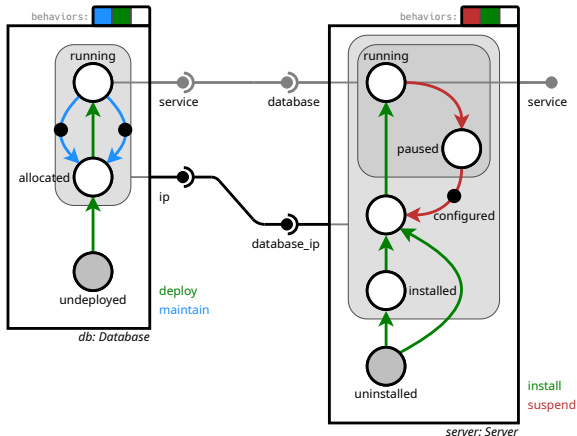
Inputs:

- reconfiguration program
- time estimations for transitions

Output: execution time prediction

Solution:

1. Dependency graph Generalization
2. Critical path computation (longest path)

*Deployment of a minimal OpenStack: 11 components, 36 services in total*



- Results on three nodes Ecotype (Nantes) of Grid'5000

- Comparison to Kolla-Ansible (production tool), and Aeolus (literature)

- Reproducible experiments on Grid'5000

# Evaluation on the deployment of OpenStack



[concerto]

[ansible]

[aeolus]

*Real use-case extracted from the OpenStack Summit 2018*

**Initial state:** centralized MariaDB running



### decentralization

- replaces centralized DB with a Galera cluster
- requires a backup of the data, and a restart of the master node

Results on nodes of UvB (Sophia) of Grid'5000 ($2\times6$-core Intel Xeon X5670 CPUs, 96 GB RAM, 250 GB HDD, internal 40 Gbps InfiniBand, external 1 Gbps)
Reproducible experiments

# Model checking on a Concerto deployment program

## Verification of deployments

Hypothesis: *the deployment of a distributed software system is written in CONCERTO and the developer wants to verify its safety and enhance its efficiency*



1. how to verify some safety properties on the deployment assembly?

2. how to enhance the efficiency without running the deployment many times?

Goal: Study the use of model checking to help solving the challenges.

**Qualitative properties**

- `deployability` (inevitability)
- `sequentiality` (observer subnet)
- `forbidden` (observer subnet)

**Quantitative properties**

- `parallelism`
  $\max(\sum(\text{reachable markings}))$
- `gantt` boundaries
  min/max costs + causality

5 versions of the OpenStack deployment successively enhanced with the tool



| $v_0$ | $v_1$ | $v_2$ | $v_3$ | $v_4$ |
|---|---|---|---|---|
| deadlock | nova improvement | nova wrong dep | mbd improvement | neutron improvement |



(a) 1-naive with critical path: *nova deploy, nova register, kst deploy, mariadb deploy, haproxy deploy*

(b) 2-nova with critical path: *nova deploy, nova upg-db, nova register, kst deploy, mariadb deploy, haproxy deploy*

Experiments conducted with the model checker Romeo

| | 0-deadlock | 1-naive | 2-nova | 3-nova | 4-nova-mdb |
|---|---|---|---|---|---|
| CONCERTO places | 27 | 27 | 28 | 28 | 29 |
| CONCERTO transitions | 22 | 22 | 25 | 25 | 28 |
| CONCERTO connections | 30 | 30 | 30 | 30 | 30 |
| Petri net places | 113 | 113 | 124 | 124 | 134 |
| Petri net transitions | 75 | 75 | 84 | 84 | 92 |
| Transformation time (ms) | 1.6 | 1.6 | 1.8 | 1.7 | 1.5 |
| Deployability | False | True | True | True | True |
| Resolution time (s) | 0 | 41.6 | 78.7 | 88.7 | 152.6 |
| Parallelism nova | - | 1 | 2 | 2 | 2 |
| Resolution time (s) | - | 42.1 | 82.7 | 93.6 | 154.3 |
| Parallelism full | - | 10 | 11 | 11 | 12 |
| Resolution time (s) | - | 43.2 | 86.1 | 98.4 | 162.9 |
| Gantt & critical path | - | Fig | Fig | Fig | Fig |
| Resolution time (s) | - | 130.1 | 266.9 | 275.4 | 588.1 |
| Boundaries | - | [575,615] | [518,554] | [400,423] | [377,398] |
| Resolution time (s) | - | 130.1, 128.8 | 266.9, 269.7 | 275.4, 267.6 | 588.1, 580.8 |

## Generalization to Concerto?

**Generalization of the approach to Concerto?**

- difficult to generalize to any reconfiguration program because of the state explosion when handling the combination of multiple behaviors per components;
- other kinds of reconfiguration patterns may be studied such as scaling, rolling upgrade etc., but may not offer an acceptable solving time.

**Synthesis**

Instead of verifying a CONCERTO reconfiguration program, we synthesize a correct reconfiguration program.

# Synthesis of Concerto reconfiguration programs

## Reconfiguration synthesis

*Work of Simon Robillard, postdoc (associate professor in Montpellier starting in September)*

**Reconfiguration synthesis**

- **input:** a reconfiguration goal
  - set of behaviors to execute on designated components
  - constraints on the final state of ports
- **output:** a correct reconfiguration script
  - `pushB` requests
  - `wait` commands
- out of scope: component creations/deletions/(dis)connections
  - usual for safety reasons to handle topological changes before and after behaviors requests and synchronization (e.g., deployment)

Solution: a three-phases algorithm

### Reconfiguration goal

1. dep1 and dep2 must execute update

2. all ports should be active at the end of the reconfiguration

**incomplete list of behaviors, a partial specification, is enough**

- find a sequence of behaviors that satisfies the goal

- enumerate and analyze possible sequences

- no solution $\implies$ failure of the procedure

- multiple solutions $\implies$ pick one according to some optimization criterion

  - shortest sequence
  - shortest estimated execution time
  - fewest port requirements



*dep1*

**solution for this example**
[update, deploy]

## 2. Correct global schedule of behaviors

- **Goal**: find a global schedule
- **assumption**: reconfiguration plan with *steps*
  - at most 1 behavior per component/step
  - each step followed by global synchronization
- **problem: assign for each required behavior a given step to schedule it**

```
1  pushB ( server , suspend )
2  waitAll ()
3  pushB ( dep1 , update )
4  pushB ( dep2 , update )
5  waitAll ()
6  pushB ( dep1 , deploy )
7  pushB ( dep2 , deploy )
8  waitAll ()
9  pushB ( server , deploy )
10 waitAll ()
```

What is expected in the second phase with 3 steps.

### SMT encoding

constraints extracted from the internal nets on sequentiality of behaviors, port requirements and status, incompatible behaviors etc.

*Missing behaviors and missing port requirements are detected during this phase*

## 3. Correct reduction of the number of synchronizations

*Solution get from phase 2*

```
1  pushB ( server , suspend )
2  waitAll ()
3  pushB ( dep1 , update )
4  pushB ( dep2 , update )
5  waitAll ()
6  pushB ( dep1 , deploy )
7  pushB ( dep2 , deploy )
8  waitAll ()
9  pushB ( server , deploy )
10 waitAll ()
```

*Final solution expected from phase 3*

```
1  pushB ( server , suspend )
2  pushB ( dep1 , update )
3  pushB ( dep2 , update )
4  pushB ( dep1 , deploy )
5  pushB ( dep2 , deploy )
6  wait ( dep1 )
7  wait ( dep2 )
8  pushB ( server , deploy )
9  wait ( server )
```

1. replace global synchronization barriers by targeted ones
2. delay barriers whenever possible by following correct rules

**Reconfiguration scenario**
update database & restore system to
working state

**Result**

- correct reconfiguration script generated in 2.49 seconds
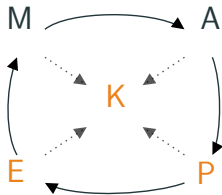- contains only synchronization barriers that are needed
- 12 behaviors on 5 components

# Conclusion and perspectives

## Conclusion

- Complexity of the reconfiguration coordination problem
- Goals: efficiency, safety
- CONCERTO reconfiguration model
- Model-checking on CONCERTO deployments
- CONCERTO reconfiguration program synthesis
- Running use-case on OpenStack

## Perspectives

- Reconfiguration patterns to improve the scalability of verification and synthesis approaches
  - deployment, scaling, rolling upgrade, substitution etc.,
  - to combine with component patterns (e.g., Docker containerized component).
- Other aspects of self-adaptation in a safe and efficient manner

# Thank you!

[1] *Predictable Efficiency for Reconfiguration of Service-Oriented Systems with Concerto*. Maverick Chardet, Hélène Coullon, Christian Perez. In CCGrid 2020.

[2] *Toward Safe and Efficient Reconfiguration with Concerto*. Maverick Chardet, Hélène Coullon, Simon Robillard. In journal SCP, 2020.

[3] *Enhancing Separation of Concerns, Parallelism, and Formalism in Distributed Software Deployment with Madeus*. Maverick Chardet, Hélène Coullon, Christian Perez, Dimitri Pertin, Charlène Servantie, Simon Robillard. [preprint]

[4] *Integrated Model-checking for the Design of Safe and Efficient Distributed Software Commissioning*. Hélène Coullon, Didier Lime, Claude Jard. In iFM 2019, Bergen, Norway.

[5] *Madeus: A formal deployment model*. Maverick Chardet, Hélène Coullon, Christian Perez and Dimitri Pertin. In 4PAD 2018 (hosted at HPCS 2018).

# Backup

## Related work

### Lifecycle

- fixed lifecycle: TOSCA, DEPLOYWARE, SMARTFROG, ENGAGE
  - easier to use, less flexible
- programmable lifecycle: AEOLUS, ANSIBLE (DevOps configuration tool)
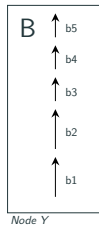  - more difficult to use, more flexible

### Dependencies

1. same component level: ANSIBLE
2. component level: TOSCA, DEPLOYWARE (DB $\ll$ WS)
3. lifecycle level: TOSCA, ENGAGE, AEOLUS (DB(3) $\ll$ WS(4))
4. intra-lifecycle level: CONCERTO

**level1: multiple nodes, same action**

- no dependencies declared
- procedural execution order
- parallelism for the same component
- ANSIBLE

### level2: level1+non-dependent components

- dependencies at the component level ($A \ll B$)
- DEPLOYWARE, (basic) TOSCA, ENGAGE

## Performance through parallelism and dependencies

### level3: level1 + level2 + inter-component

- dependencies at the lifecycle level between components ($a3 \ll b3$)
- (advanced) TOSCA, AEOLUS

**level 4: level1 + level2 + level3 + intra-component**

- parallelism within the lifecycle of one component (b1 $\ll$ b2)
- CONCERTO



The finer the dependencies granularity is, the better is the efficiency

# Control components in practice

Written by the component developers

```python
class Server(Component):
    def create(self):
        self.places = ['uninstalled','installed','configured','running','paused']

        self.initial_place = 'uninstalled'

        self.behaviors = ['b_install','b_suspend']

        self.transitions = {
            'install1': ('uninstalled','installed','b_install',self.install1),
            'install2': ('uninstalled','configured','b_install',self.install2),
            'configure': ('installed','configured','b_install',self.configure),
            'start': ('configured','running','b_install',self.start),
            'suspend1': ('running','paused','b_suspend',self.suspend1),
            'suspend2': ('paused','configured','b_suspend', self.suspend2)
        }
```

# Control components in practice

Written by the component developers

```python
1  class Server(Component):
2      def create(self):
3          ...
4
5          self.dependencies = {
6              'database_ip': (DepType.USE, ['installed','configured','running','paused']),
7              'database': (DepType.USE, ['running','paused']),
8              'service': (DepType.PROVIDE, ['running'])
9          }
10
11     # Definition of the actions
12     def install1(self):
13         remote = SSHClient()
14         remote.connect(host, user, pwd)
15         remote.exec_command(cmd)
16         ...
```

## Deployments pattern with Concerto

### Assumptions

- for each component there is one behavior that leads from an uninstalled place to a running place, namely `deploy`;

```
1 for i in [1,n]
2     add(i)
```

- the final assembly is specified by the user;

```
1 connect as specified by the user
```

### Semantics

Simply apply the Concerto semantics on a single behavior per component

```
1 for i in [1,n]
2     pushB(i, deploy)
3 waitall
```

HALP automatically transformed to TCTL (Time Computational Tree Logic) formulae

**Qualitative properties**

- deployability $\longrightarrow$ inevitability
- sequentiality $\longrightarrow$ observer subnet + invariant
- forbidden $\longrightarrow$ observer subnet

**Quantitative properties**

- parallelism $\longrightarrow$ max($\sum$(reachable markings))
- gantt boundaries: min/max costs + causality in the trace to get the critical path

Constraints added to the problem:

**sequentiality of behaviors**
from step 1. $int(schedule(dep1.update)) < int(schedule(dep1.deploy))$

**port requirements at the beginning of behaviors**
$\neg active_u(schedule(dep1.update))$

**separation of behaviors with incompatible port effects**
$schedule(server.deploy) \neq schedule(dep1.update)$

**port status after behaviors**
$active_p(succ(schedule(dep1.deploy)))$

## Missing port requirements in 2.

**Problem:** some unsatisfied ports requirements may make scheduling impossible

### Example

- we determined that dep1 and dep2 should execute [update, deploy]
- but the updates can't be executed while the server is relying on the provide ports

### Solution:

1. deduce new individual component goals
   - go to state that satisfies missing port requirement
   - go to state that satisfies final port constraints
2. extend the set of behaviors to schedule and go back to phase 1