

XPFlow (Experimental workflow)

Tomasz Buchert



We all know how frustrating experimenting can be.



That's because experiments in distributed systems are:

- time-consuming
- difficult to do correctly
- complex and incomprehensible
- failure-prone

With tools like Chef and Puppet:

- a human factor is nearly removed
- systems are built from modules
- the configuration is reproducible

But reproducibility does not necessarily imply **descriptiveness**.
It does not imply **ease of understanding** either.

Many tools to manage experiments exist:

- Expo
- g5k-campaign
- OMF
- Plush
- ... among many others

They are based on different paradigms.

Bottom-up vs top-down approach

Most of these tools use **bottom-up design**.

What about a **top-down** approach?

- ① Start with high-level description of the experiment.
- ② Implement low-level details.
- ③ Run the experiment.
- ④ Improve if necessary and reiterate.

There already exists an approach like this.

Business Process Management

Business Process Management is about:

- understanding an organization
- modeling its processes as **workflows**
- **executing** processes and **monitoring** them
- **improving** organizational **activities**
- redesigning **processes** to make them:
 - cheaper
 - faster
 - less defective

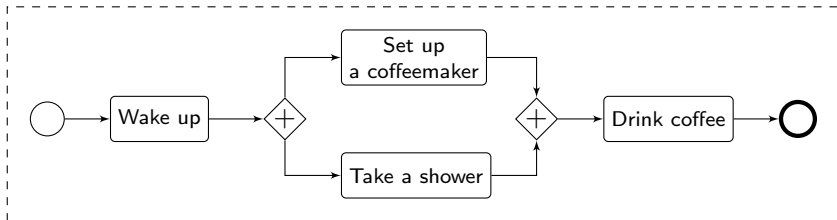


In this talk I will present **XPFlow**:

- a new experimentation engine
- based on **Business Process Modeling and Management**

There are 2 main concepts in XPFlow:

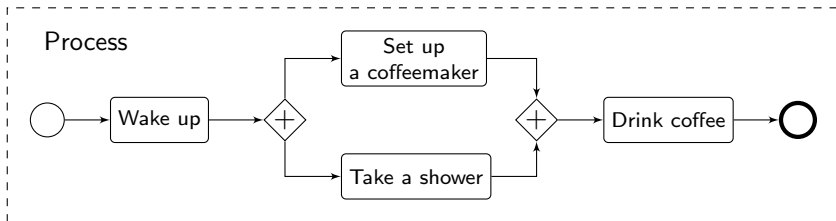
- Processes – high-level description of an experiment:
 - workflows written in a DSL
 - orchestrate other processes and activities
- Activities – low-level building blocks of experiments:
 - do real hard work
 - written in Ruby



XPFlow concepts

There are 2 main concepts in XPFlow:

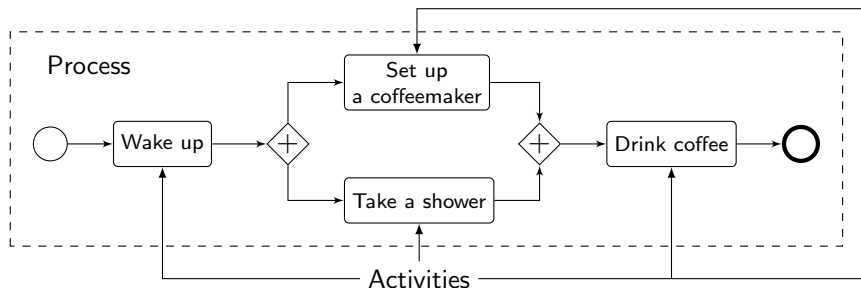
- Processes – high-level description of an experiment:
 - workflows written in a DSL
 - orchestrate other processes and activities
- Activities – low-level building blocks of experiments:
 - do real hard work
 - written in Ruby



XPFlow concepts

There are 2 main concepts in XPFlow:

- Processes – high-level description of an experiment:
 - workflows written in a DSL
 - orchestrate other processes and activities
- Activities – low-level building blocks of experiments:
 - do real hard work
 - written in Ruby

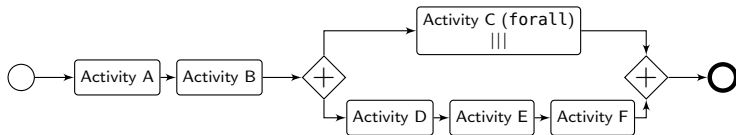


The DSL for processes features different **workflow patterns**:

- running activities and other processes (run),
- running activities in order or in parallel (sequence, parallel),
- conditional expressions (if, switch)
- running sequential and parallel loops (loop, foreach, forall),
- error handling (try, checkpoint).

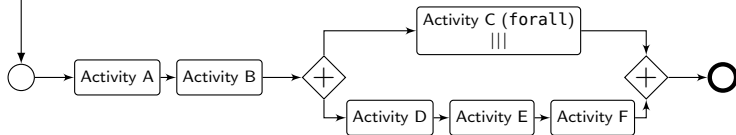
Some of them are taken directly from BPM.

Workflow patterns (example)

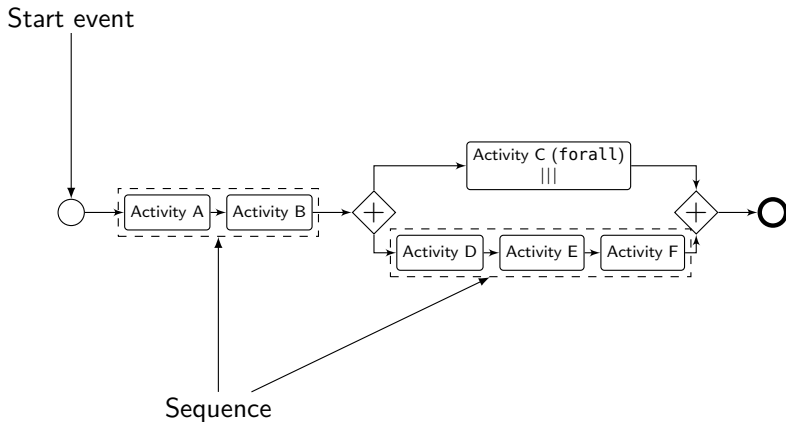


Workflow patterns (example)

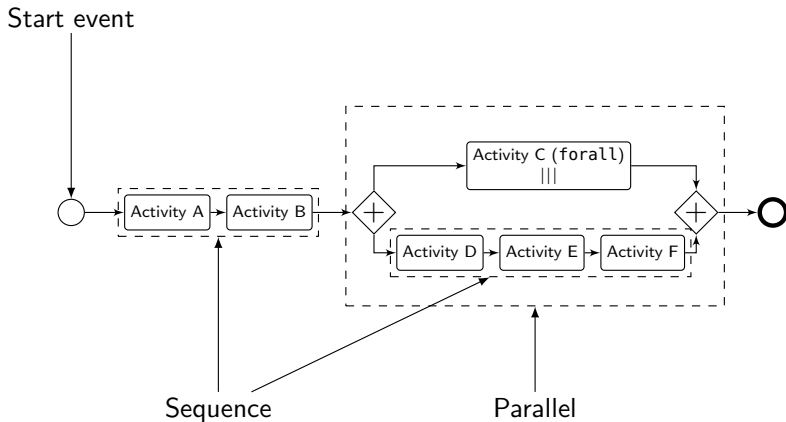
Start event



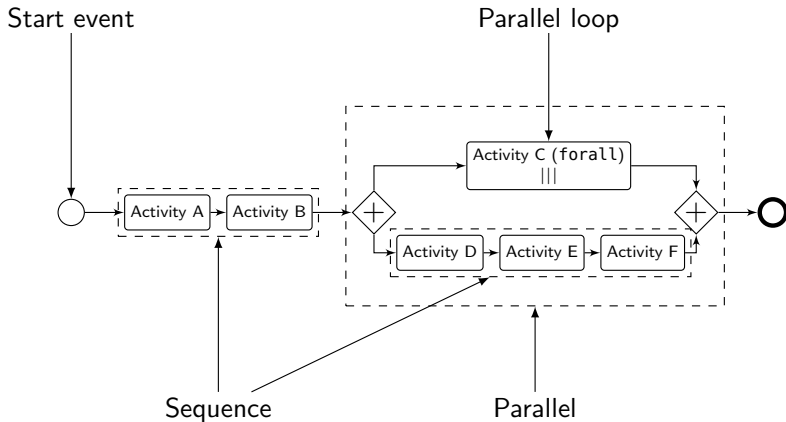
Workflow patterns (example)



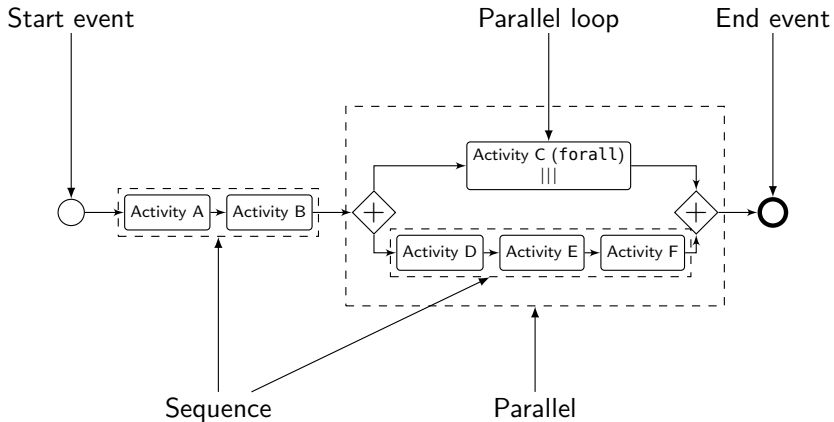
Workflow patterns (example)



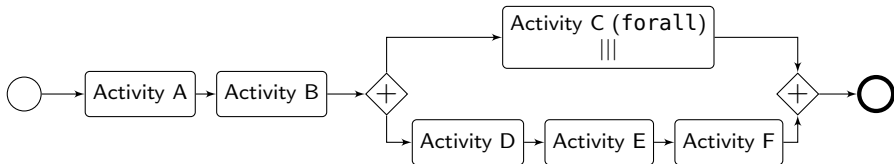
Workflow patterns (example)



Workflow patterns (example)

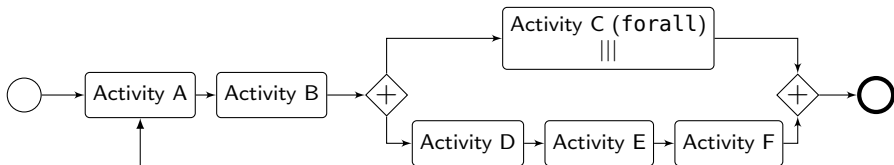


Workflow patterns (example, cont.)



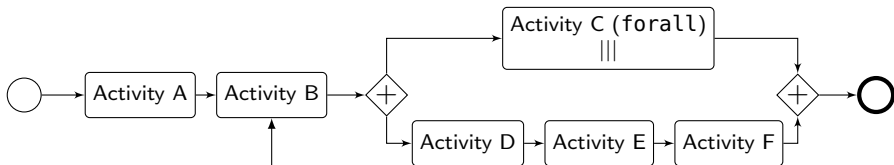
```
process :workflow do |array|
  run :a
  run :b
  parallel do
    forall array do |x|
      run :c, x
    end
    sequence do
      run :d
      run :e
      run :f
    end
  end
end
end
```

Workflow patterns (example, cont.)



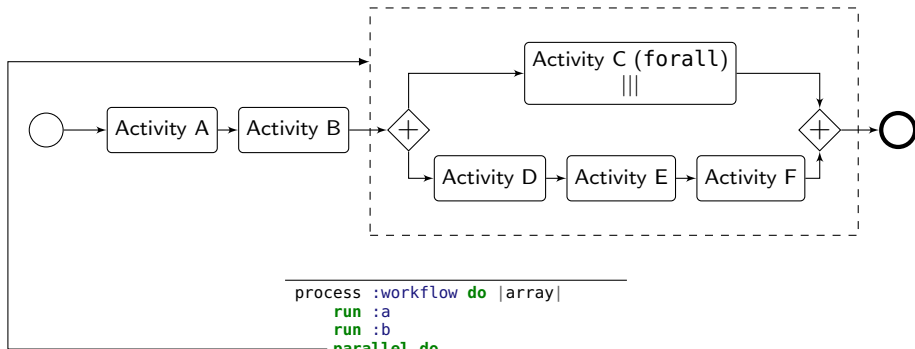
```
process :workflow do |array|  
  run :a  
  run :b  
  parallel do  
    forall array do |x|  
      run :c, x  
    end  
    sequence do  
      run :d  
      run :e  
      run :f  
    end  
  end  
end  
end
```

Workflow patterns (example, cont.)



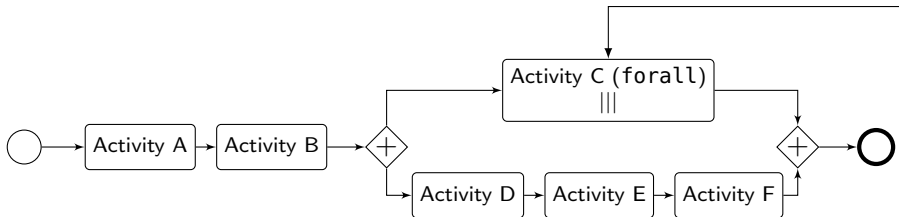
```
process :workflow do |array|  
  run :a  
  run :b  
  parallel do  
    forall array do |x|  
      run :c, x  
    end  
    sequence do  
      run :d  
      run :e  
      run :f  
    end  
  end  
end  
end
```

Workflow patterns (example, cont.)



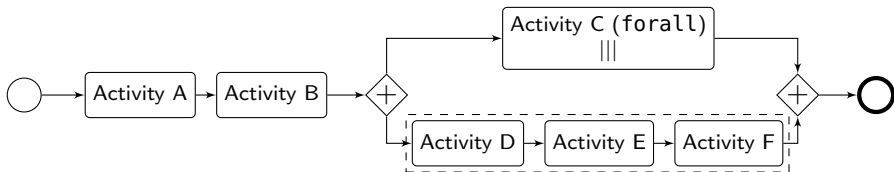
```
process :workflow do |array|  
  run :a  
  run :b  
  parallel do  
    forall array do |x|  
      run :c, x  
    end  
    sequence do  
      run :d  
      run :e  
      run :f  
    end  
  end  
end  
end
```

Workflow patterns (example, cont.)



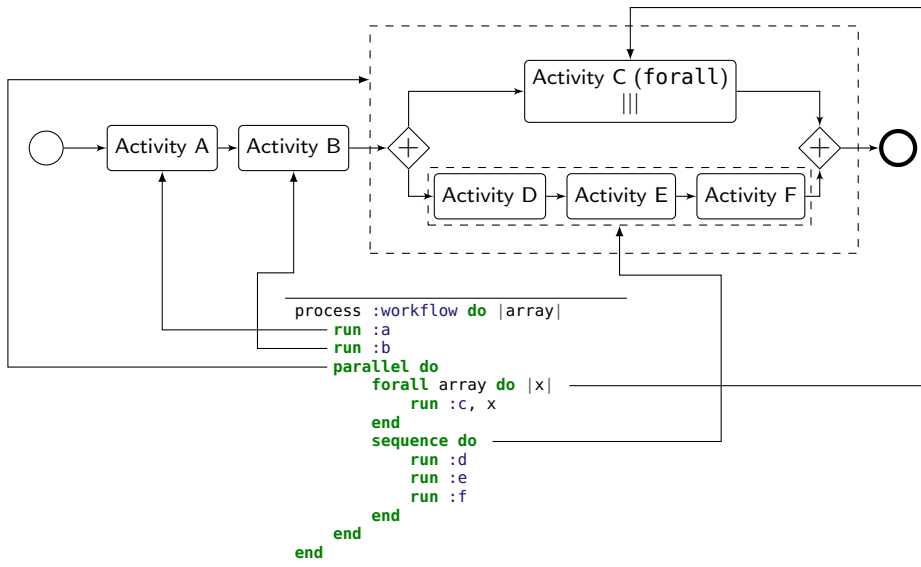
```
process :workflow do |array|  
  run :a  
  run :b  
  parallel do  
    forall array do |x|  
      run :c, x  
    end  
    sequence do  
      run :d  
      run :e  
      run :f  
    end  
  end  
end  
end
```

Workflow patterns (example, cont.)



```
process :workflow do |array|  
  run :a  
  run :b  
  parallel do  
    forall array do |x|  
      run :c, x  
    end  
    sequence do  
      run :d  
      run :e  
      run :f  
    end  
  end  
end  
end
```

Workflow patterns (example, cont.)



Minimal Grid'5000 example

```
#!/usr/bin/env xpfloor

use :g5k

process :clean_all do
  sites = run :g5k.sites
  forall sites do |s|
    jobs = run :g5k.my_jobs, s
    log "Cleaning ", s, " from ", (length_of jobs), " jobs..."
    run :g5k.release_all, s
  end
end

main :clean_all
```

Assumes that xpfloor is in your \$PATH.

Slightly less minimal Grid'5000 example

```
#!/usr/bin/env xpfLOW

use :g5k

process :build_kernel do |sources|
  run :ensure_exists, sources
  r = run :g5k.reserve_nodes,
    :nodes => 1,
    :time => '1h',
    :site => 'nancy'
  nodes = run :g5k.nodes, r
  node = (first_of nodes)
  checkpoint :reserved
  log 'Our node is: ', node
  run :copy_sources,
    sources, node
  checkpoint :source_copied
  run :make_kernel, node
  checkpoint :kernel_made
  run :copy_debs, node
  run :g5k.release, r
end

main :build_kernel, :str
```

```
activity :copy_sources do |sources, node|
  run 'g5k.copy', sources,
    node, '/tmp/linux.tar.gz'
  run 'g5k.bash', node do
    cd '/tmp'
    run 'rm -rf linux-build'
    make_dirs 'linux-build'
    untar 'linux.tar.gz',
      'linux-build'
  end
end

activity :make_kernel do |node|
  run 'g5k.bash', node do
    cd '/tmp/linux-build'
    cd dirs.first
    run 'make deb-pkg -j 5'
  end
end

activity :copy_debs do |node|
  run 'g5k.retrieve', node,
    '/tmp/linux-build/*.deb'
  run 'g5k.bash', node do
    run 'rm -rf /tmp/linux-build'
  end
end
```

XPFlow gives some means to cope with failures:

- snapshotting:
 - saves a state of an experiment for future use
 - shortens a development's cycle
- retry policy:
 - retries a failed subprocess execution
 - improves reliability

```
process :snapshotting do
  run :long_deployment
  checkpoint :d
  run :experiment
end
```

```
process :retrying do
  try :retry => 5 do
    run :tricky_activity
  end
end
```

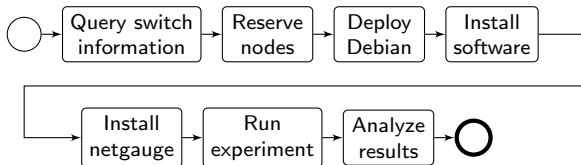
Both features require **idempotency**.

Example of an experiment

Measure the *effective bisection bandwidth* of a switch.

- 1 Get names of all nodes connected to the switch.
- 2 Reserve the nodes.
- 3 Deploy Debian OS.
- 4 Install necessary software.
- 5 Compile and install *netgauge*.
- 6 Run the experiment.
- 7 Analyze results.

An experiment workflow

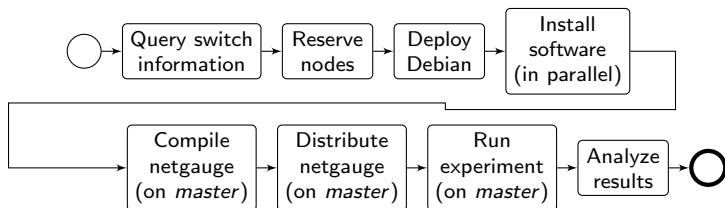


Few notes:

- each node must have some software installed
- each node must have *netgauge* installed ...
- ... but one node is enough to compile it
- one node must launch MPI application

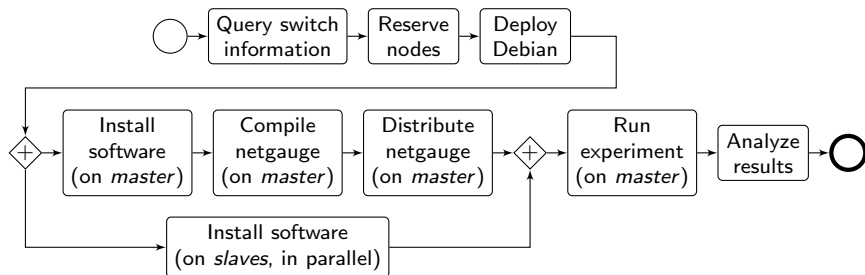
We will introduce a *master* node and *slave* nodes.

An experiment workflow



Another observation: compilation can run in parallel with installation of software on the *slave* nodes.

An experiment workflow



This workflow describes our experiment.

The last thing to do is to express that in XPFlow.

An experiment workflow - DSL representation

```
process :exp do |site, switch|
  s = run g5k.switch, site, switch
  ns = run g5k.nodes, s
  r = run g5k.reserve_nodes,
        :nodes => ns, :time => '2h',
        :site => site, :type => :deploy
  master = (first_of ns)
  rest = (tail_of ns)
  run g5k.deploy,
        r, :env => 'squeeze-x64-nfs'
  checkpoint :deployed
  parallel :retry => true do
    forall rest do |slave|
      run :install_pkgs, slave
    end
    sequence do
      run :install_pkgs, master
      run :build_netgauge, master
      run :dist_netgauge,
          master, rest
    end
  end
  checkpoint :prepared
  output = run :netgauge, master, ns
  checkpoint :finished
  run :analysis, output, switch
end
```

An experiment workflow - DSL representation

```
process :exp do |site, switch|
  s = run g5k.switch, site, switch
  ns = run g5k.nodes, s
  r = run g5k.reserve_nodes,
      :nodes => ns, :time => '2h',
      :site => site, :type => :deploy
  master = (first_of ns)
  rest = (tail_of ns)
  run g5k.deploy,
      r, :env => 'squeeze-x64-nfs'
  checkpoint :deployed
  parallel :retry => true do
    forall rest do |slave|
      run :install_pkgs, slave
    end
    sequence do
      run :install_pkgs, master
      run :build_netgauge, master
      run :dist_netgauge,
          master, rest
    end
  end
  checkpoint :prepared
  output = run :netgauge, master, ns
  checkpoint :finished
  run :analysis, output, switch
end
```

Activity :install_pkgs

```
activity :install_pkgs do |node|
  log 'Installing packages on ', node
  run 'g5k.bash', node do
    aptget :update
    aptget :upgrade
    aptget :purge, 'mx'
  end
end
```

An experiment workflow - DSL representation

```
process :exp do |site, switch|
  s = run g5k.switch, site, switch
  ns = run g5k.nodes, s
  r = run g5k.reserve_nodes,
        :nodes => ns, :time => '2h',
        :site => site, :type => :deploy
  master = (first_of ns)
  rest = (tail_of ns)
  run g5k.deploy,
        r, :env => 'squeeze-x64-nfs'
  checkpoint :deployed
  parallel :retry => true do
    forall rest do |slave|
      run :install_pkgs, slave
    end
    sequence do
      run :install_pkgs, master
      run :build_netgauge, master
      run :dist_netgauge,
          master, rest
    end
  end
  checkpoint :prepared
  output = run :netgauge, master, ns
  checkpoint :finished
  run :analysis, output, switch
end
```

Activity :build_netgauge

```
activity :build_netgauge do |master|
  log "Building netgauge on #{master}"
  run 'g5k.copy', NETGAUGE, master, '~',
  run 'g5k.bash', master do
    build_tarball NETGAUGE, PATH
  end
  log "Build finished."
end
```

An experiment workflow - DSL representation

```
process :exp do |site, switch|
  s = run g5k.switch, site, switch
  ns = run g5k.nodes, s
  r = run g5k.reserve_nodes,
      :nodes => ns, :time => '2h',
      :site => site, :type => :deploy
  master = (first_of ns)
  rest = (tail_of ns)
  run g5k.deploy,
      r, :env => 'squeeze-x64-nfs'
  checkpoint :deployed
  parallel :retry => true do
    forall rest do |slave|
      run :install_pkgs, slave
    end
    sequence do
      run :install_pkgs, master
      run :build_netgauge, master
      run :dist_netgauge,
          master, rest
    end
  end
  checkpoint :prepared
  output = run :netgauge, master, ns
  checkpoint :finished
  run :analysis, output, switch
end
```

Activity :dist_netgauge

```
activity :dist_netgauge do |m, s|
  master, slaves = m, s
  run 'g5k.dist_keys', master, slaves
  run 'g5k.bash', master do
    distribute BINARY,
        DEST, 'localhost', slaves
  end
end
```

An experiment workflow - DSL representation

```
process :exp do |site, switch|
  s = run g5k.switch, site, switch
  ns = run g5k.nodes, s
  r = run g5k.reserve_nodes,
        :nodes => ns, :time => '2h',
        :site => site, :type => :deploy
  master = (first_of ns)
  rest = (tail_of ns)
  run g5k.deploy,
    r, :env => 'squeeze-x64-nfs'
  checkpoint :deployed
  parallel :retry => true do
    forall rest do |slave|
      run :install_pkgs, slave
    end
    sequence do
      run :install_pkgs, master
      run :build_netgauge, master
      run :dist_netgauge,
        master, rest
    end
  end
  checkpoint :prepared
  output = run :netgauge, master, ns
  checkpoint :finished
  run :analysis, output, switch
end
```

Activity :netgauge

```
activity :netgauge do |master, nodes|
  log "Running experiment..."
  out = run 'g5k.bash', master do
    cd PATH
    mpirun nodes, "./netgauge"
  end
  log "Experiment done."
end
```

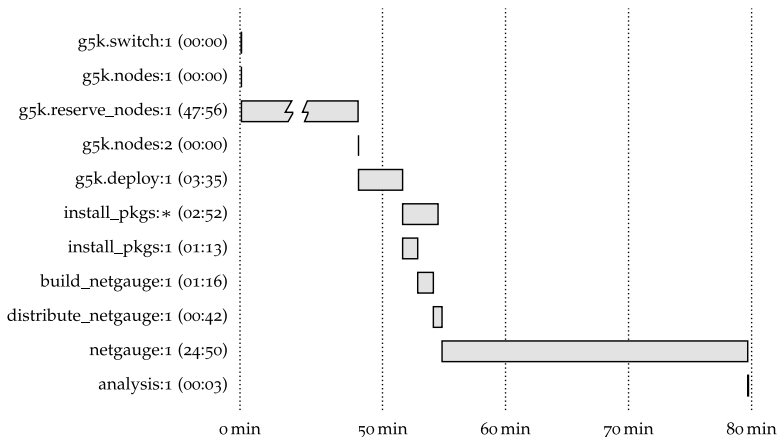
Running the experiment

The experiment runs on Grid'5000 frontend or on your local machine.

```
[ 11:15:52.940 ] Started activity g5k.switch:1.  
[ 11:15:53.418 ] Finished activity g5k.switch:1 (0.478 s).  
[ 11:15:53.419 ] Process exp: Experimenting with switch: sgraphene2  
[ 11:15:53.419 ] Started activity g5k.nodes:1.  
[ 11:15:53.419 ] Finished activity g5k.nodes:1 (0.000 s).  
[ 11:15:53.419 ] Started activity g5k.reserve_nodes:1.  
[ 11:15:55.837 ] Waiting for reservation 408387  
[ 11:16:02.452 ] Reservation 408387 should be available in 12 mins  
[ 11:16:02.452 ] Reservation 408387 ready  
[ 11:16:02.453 ] Finished activity g5k.reserve_nodes:1 (9.022 s).  
[ 11:16:02.453 ] Started activity g5k.nodes:2.  
[ 11:16:02.453 ] Finished activity g5k.nodes:2 (0.000 s).  
[ 11:16:02.453 ] Started activity g5k.deploy:1.  
[ 11:22:09.427 ] Finished activity g5k.deploy:1 (366.968 s).  
[ 11:22:09.429 ] Started activity install_pkgs.  
[ 11:22:09.429 ] Started activity install_pkgs:1.  
[ 11:22:09.430 ] Activity install_pkgs: Installing packages on graphene-96  
[ 11:22:09.430 ] Started activity install_pkgs:2.  
[ 11:22:09.430 ] Activity install_pkgs: Installing packages on graphene-60
```

The execution is monitored and errors reported if necessary.

Monitoring features - Gantt chart of the execution



Each activity is monitored during its execution.

Notice that `build_netgauge:1` runs in parallel with `install_pkgs:*`.

Features to implement:

- GUI to see progress and monitor
- workflow validation (using type checking)
- support for other testbed
- expressing G5K interface using workflows
- result/data provenance
- distributed workflows
- efficient operations
- improved modularity
- monitoring of the platform
- tutorial :)

In this talk I presented XPFlow.

Current features include:

- improved descriptiveness
- modularity and flexibility
- monitoring and support for common patterns
- integration with Grid'5000

In the future we will work on:

- integration with experimentation tools
- human interaction during the experiment
- efficient data broadcast and collection

Thank you for your attention. Questions?

XPFlow technical thingies:

- limitations of Ruby language (own language?)
- Ruby version hell
- access from outside G5K (also from Jenkins)
- repository at gforge (private, though)
- Jenkins at <http://ci.inria.fr> (Q: how to test?)

What about:

- treating everything as *actors* (or interfaces)?
- using functional programming to model experiments?
- having a virtual testbed?
- having an interactive mode?

Problems:

- “nested experimentation” ($G5K \rightarrow Distem \rightarrow Exp$)
- duplicated code
- many small problems
- ???

Care for demo?