

A Technical Anatomy of How OpenMPI Applications Can Inherit Fault Tolerance Using SPM.Python

Minesh B. Amin

mamin@mbasciences.com

<http://www.mbasciences.com>

PyHPC Workshop

Supercomputing Conference 2011

Seattle, Washington

Nov 18, 2011

Prologue

```
GNU/Linux [] mpirun ... ./hello_world -prefix "api"
```

Prologue

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ...

Prologue

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ... that lacks support for:

Prologue

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ... that lacks support for:

- fault tolerance

Prologue

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ... that lacks support for:

- fault tolerance
- timeout

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ... that lacks support for:

- fault tolerance
- timeout
- detection of deadlocks

Prologue

GNU/Linux [] `mpirun/hello_world -prefix "api"`



Typical OpenMPI application ... that lacks support for:

- fault tolerance
- timeout
- detection of deadlocks

⇒ Prototyping is (deeply)[∞] frustrating

Problem Statement

Prototyping should be frictionless

Problem Statement

Prototyping should be frictionless



Must use original OpenMPI application

- original source code
- original binary

Problem Statement

Prototyping should be frictionless



Must use original OpenMPI application

- original source code
- original binary



Original OpenMPI application must inherit support for:

- fault tolerance
- timeout
- detecting deadlocks

Problem Statement

Prototyping should be frictionless



Must use original OpenMPI application

- original source code
- original binary



Original OpenMPI application must inherit support for:

- fault tolerance
- timeout
- detecting deadlocks



GNU/Linux []

~~spm.python ...~~
~~mpixn ...~~

`./hello_world -prefix "api"`

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

./hello_world -prefix "api"

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpir~~in ...

./hello_world -prefix "api"

A

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

B

./hello_world -prefix "api"

A

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

./hello_world -prefix "api"

B

A

Exploiting two very different forms of parallelism:

- Using same resources
- At the same time

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

./hello_world -prefix "api"

B

A

Drop-in
replacement for
mpirun

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

./hello_world -prefix "api"

B

A

Drop-in replacement for **mpirun**

Multiple sessions of **mpirun** within a single session of **spm.python**

Problem Statement (Cont'd)

GNU/Linux []

spm.python ...
~~mpirun ...~~

./hello_world -prefix "api"

B

A

Drop-in replacement for **mpirun**

Multiple sessions of **mpirun** within a single session of **spm.python**

- Can use same resources for:
- Checkpoint based parallelism
 - What-if analysis
 - Stress testing

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the management of a collection of serial tasks which may communicate using only compatible communication primitives

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and
- the manner in which resources are obtained and released

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and
- the manner in which resources are obtained and released

serial tasks are classified in terms of either:

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the **management** of a collection of **serial tasks** which may communicate using only **compatible communication primitives**

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and
- the manner in which resources are obtained and released

serial tasks are classified in terms of either:

- **Coarse grain** ...
where tasks may not communicate prior to conclusion, or

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the management of a collection of serial tasks which may communicate using only compatible communication primitives

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and
- the manner in which resources are obtained and released

serial tasks are classified in terms of either:

- Coarse grain ...
where tasks may not communicate prior to conclusion, or
- Fine grain ...
where tasks may communicate prior to conclusion.

Terminology: "Exploiting Parallelism"

Exploiting parallelism entails the management of a collection of serial tasks which may communicate using only compatible communication primitives

management refers to policies by which:

- tasks are scheduled,
- premature terminations are handled,
- preemptive support is provided,
- communication primitives are enabled/disabled, and
- the manner in which resources are obtained and released

serial tasks are classified in terms of either:

- Coarse grain ...

Management policies codify how serial tasks are to be managed ... independent of what they may be

n, or

Terminology: "Parallel Enabling Technologies"

Means to the end

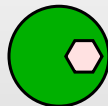
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance



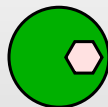
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

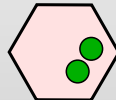
- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance



- Top-down

Hadoop Goldenorb
GraphLab

- Limited flexibility
- Fewer headaches
- Fault tolerance is inherited



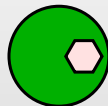
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

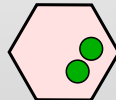
- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance



- Top-down

Hadoop Goldenorb
GraphLab

- Limited flexibility
- Fewer headaches
- Fault tolerance is inherited



- Self-contained environment

SPM.Python

- Maximum flexibility
- Fewest headaches
- Fault tolerance is inherited



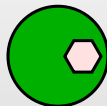
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance



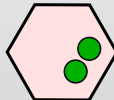
- Top-down

Hadoop Goldenorb

- Limited flexibility

N environments/installations for N Frameworks

- Fault tolerance is inherited



- Self-contained environment

SPM.Python

- Maximum flexibility
- Fewest headaches
- Fault tolerance is inherited



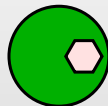
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance

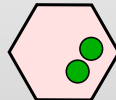


- Top-down

Hadoop Goldenorb

N environments/installations for N Frameworks

- Fault tolerance is inherited



- Self-contained environment

SPM P4H

One environment/installation, N suites of Plosures

- Fault tolerance is inherited



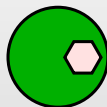
Terminology: "Parallel Enabling Technologies"

Means to the end

- Bottom-up

OpenMPI OpenMP
CUDA OpenGL

- Maximum flexibility
- Maximum headaches
- Must implement fault tolerance

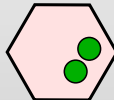


- Top-down

Hadoop Goldenorb

N environments/installations for N Frameworks

fault tolerance is inherited



- Self-contained environment

SPM

One

```
>>> createVirtualCloud -async
>>> cmdA
>>> cmdB
>>> cmdC
>>> cmdD
>>> cmdA -parallel
>>> cmdB -parallel
>>> cmdC -parallel
>>> cmdD -parallel
```



Anatomy: Timeline

GNU/Linux []

spm.python ...
~~mpfun~~

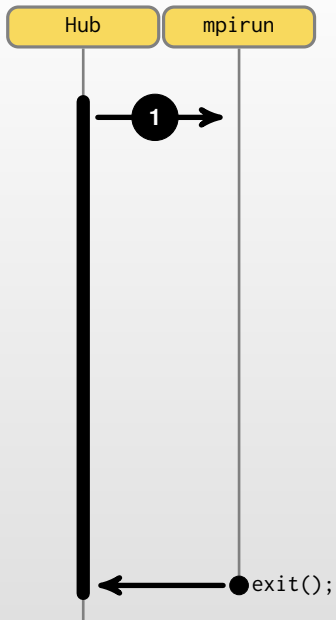
./hello_world -prefix "api"

Anatomy: Timeline (Cont'd)

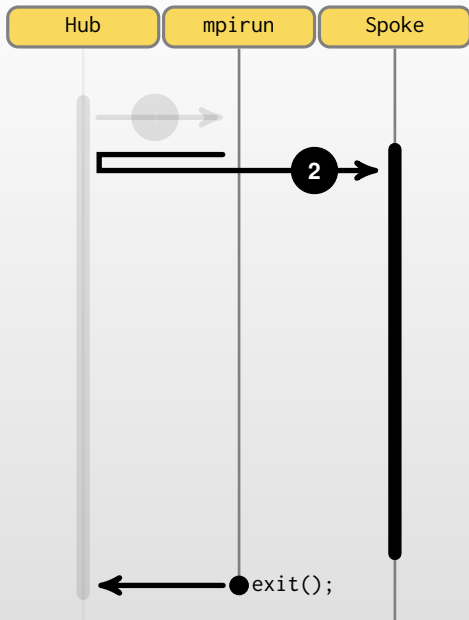
Hub



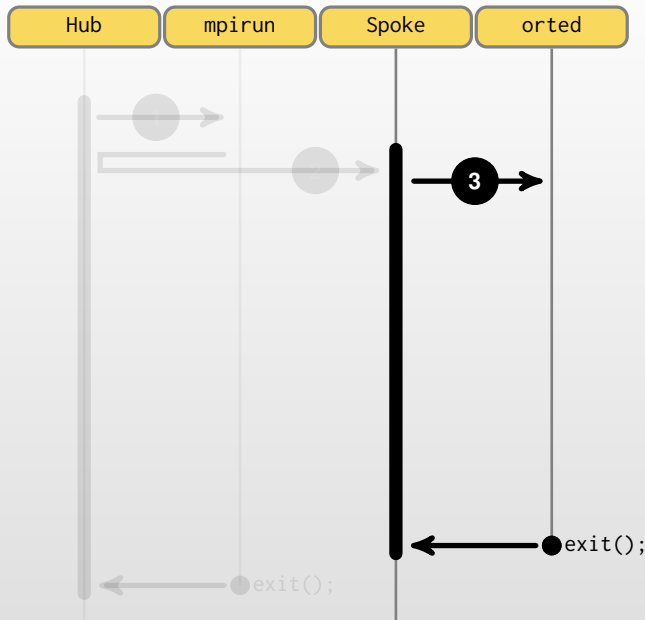
Anatomy: Timeline (Cont'd)



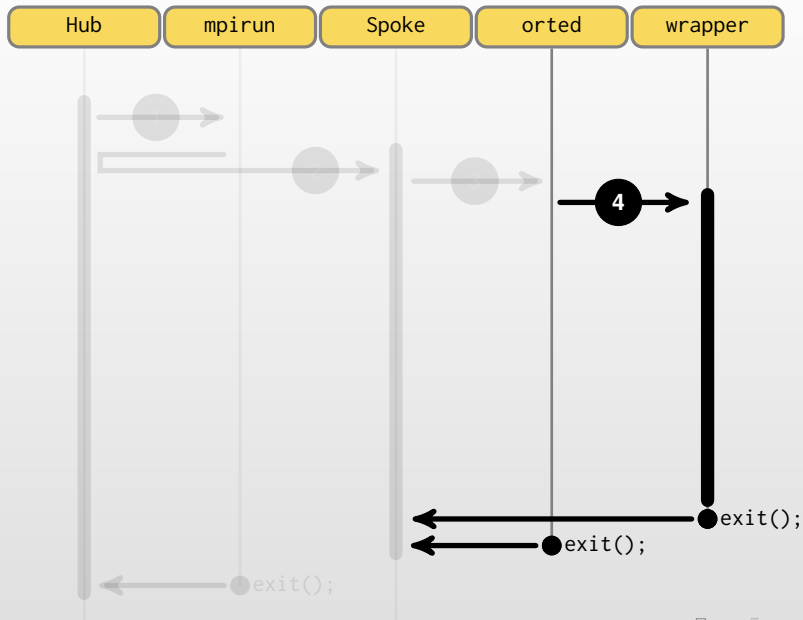
Anatomy: Timeline (Cont'd)



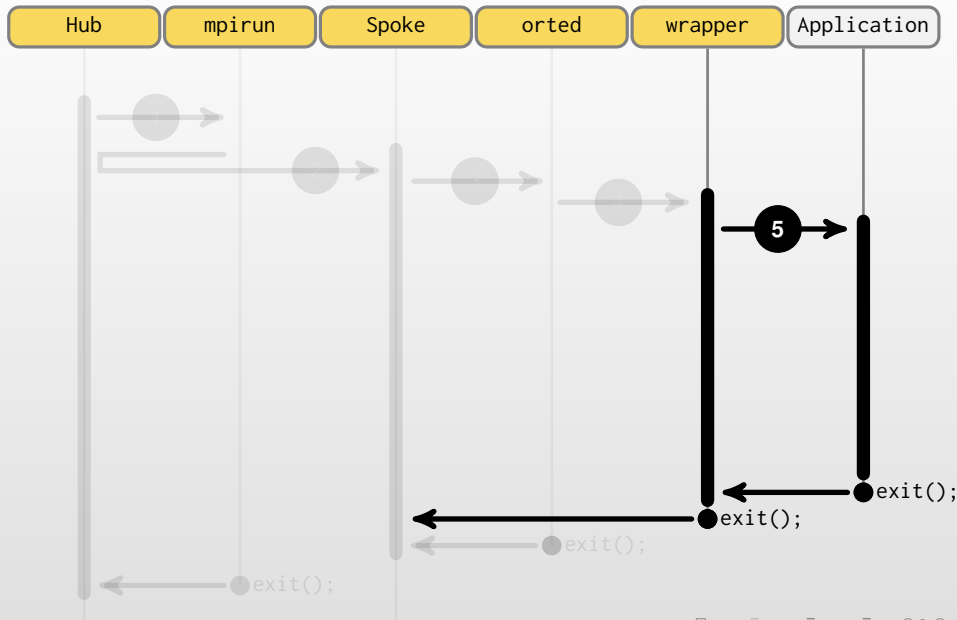
Anatomy: Timeline (Cont'd)



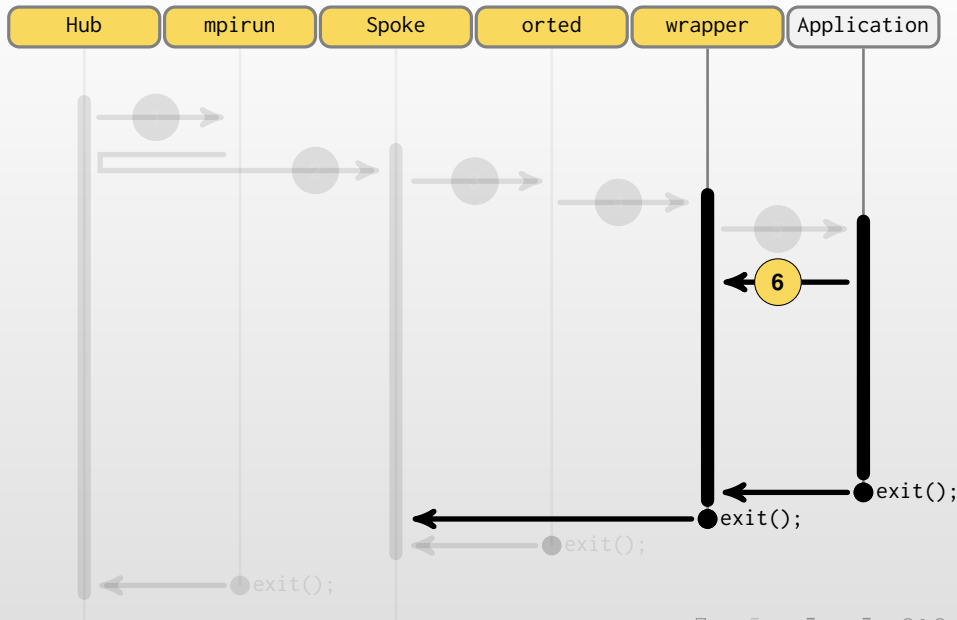
Anatomy: Timeline (Cont'd)



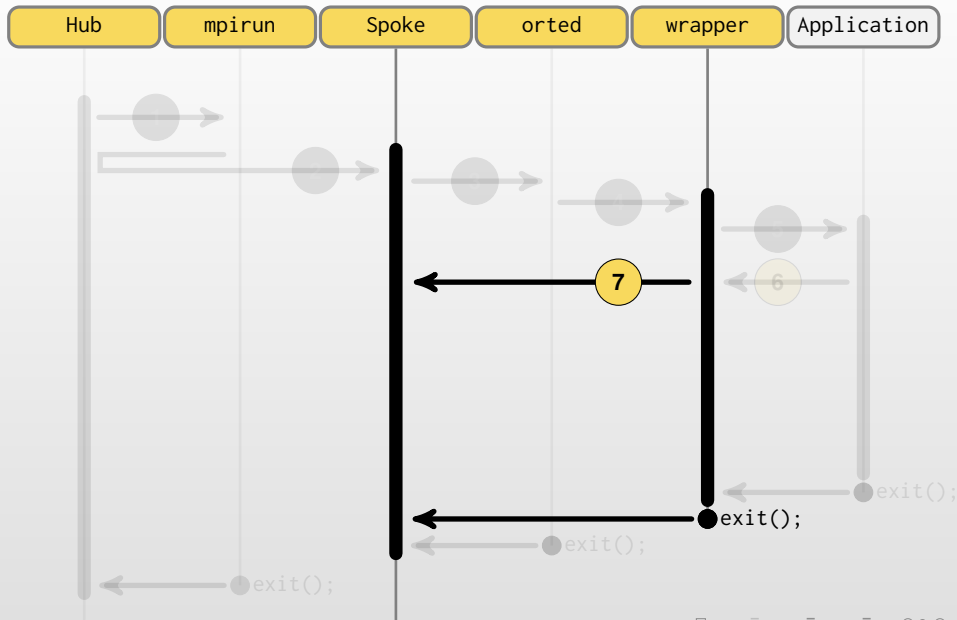
Anatomy: Timeline (Cont'd)



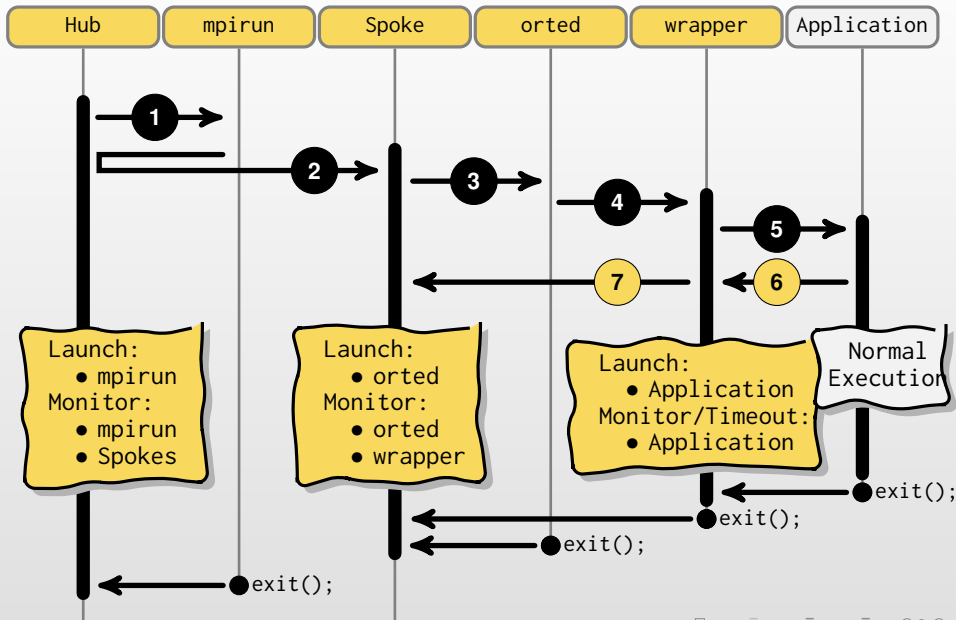
Anatomy: Timeline (Cont'd)



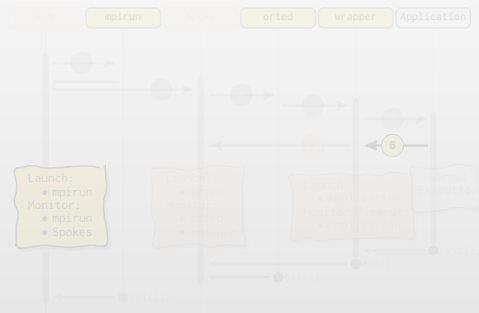
Anatomy: Timeline (Cont'd)



Anatomy: Timeline (Cont'd)



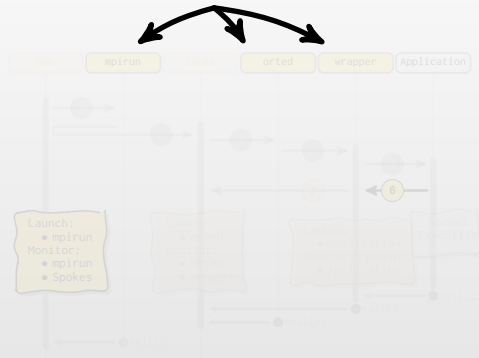
Anatomy: Breakdown



Anatomy: Breakdown

Built-in Package Management System

- Selectively change default OpenMPI env




Built-in Package Management System

```
>>> sys.path = [ ".", "@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```

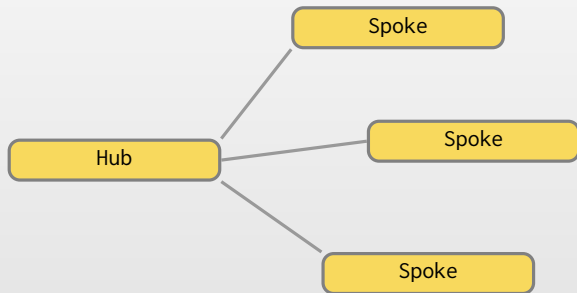
Built-in Package Management System

```
>>> sys.path = [ ".", "@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



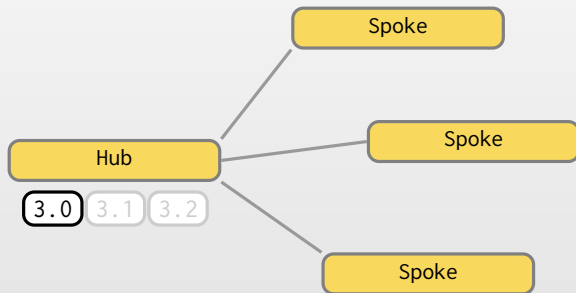
Built-in Package Management System

```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



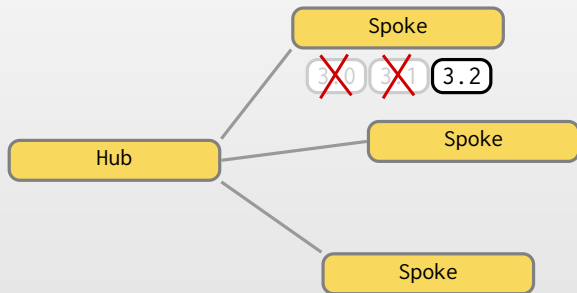
Built-in Package Management System

```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



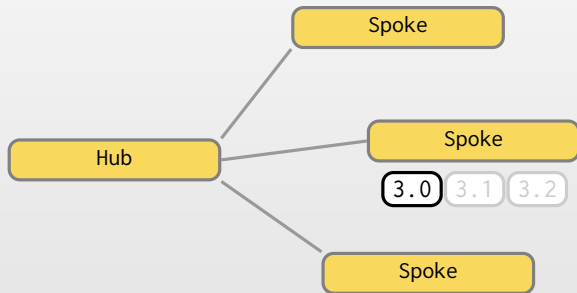
Built-in Package Management System

```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



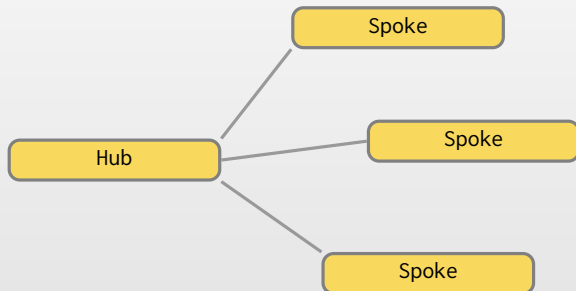
Built-in Package Management System

```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



Built-in Package Management System

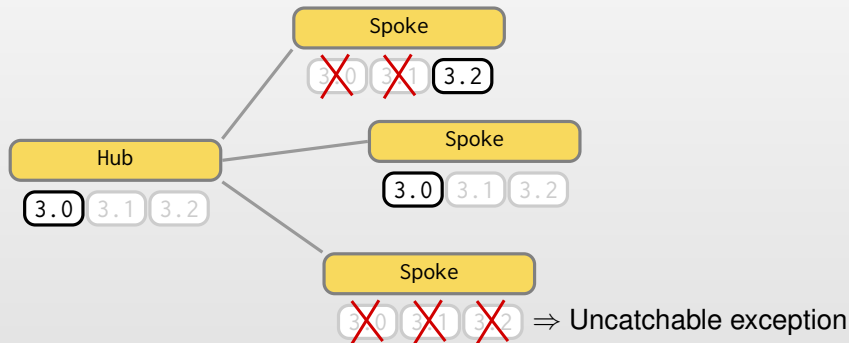
```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



~~3.0~~ ~~3.1~~ ~~3.2~~ ⇒ Uncatchable exception

Built-in Package Management System

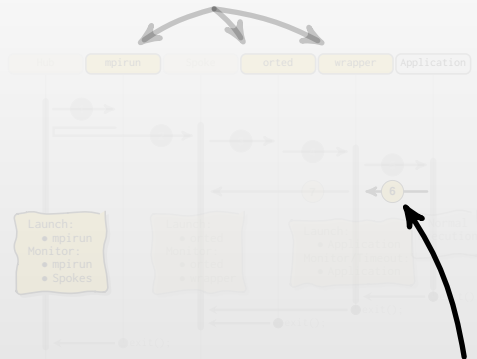
```
>>> sys.path = [ ".", "-@-/pkg.builtin", "/opt/default" ]  
>>> import pycuda
```



Anatomy: Breakdown

Built-in Package Management System

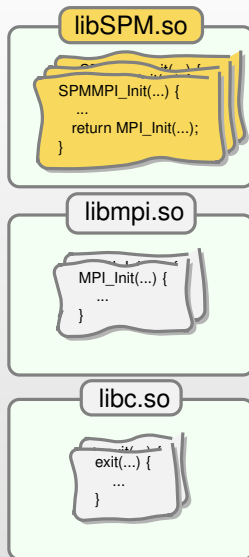
- Selectively change default OpenMPI env



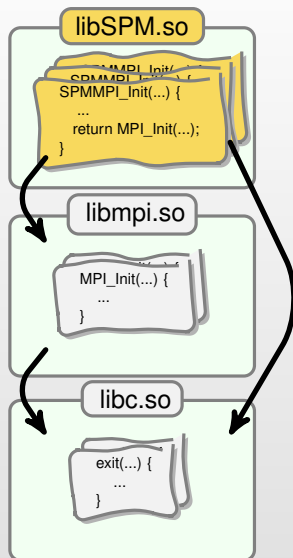
Redirection of library calls

- Augment libmpi.so, libc.so ... with libSPM.so

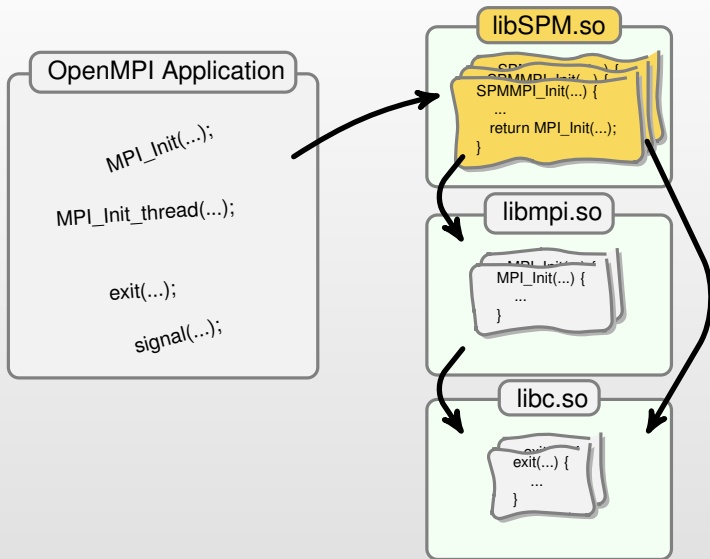
Redirecting Shared Library Calls



Redirecting Shared Library Calls



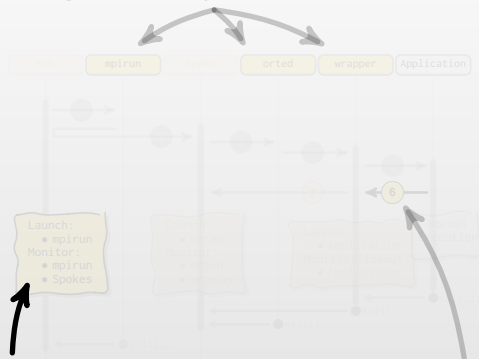
Redirecting Shared Library Calls



Anatomy: Breakdown

Built-in Package Management System

- Selectively change default OpenMPI env



Second Parallel Capability

- ~ 60-line python script
- Authored by developer

Redirection of library calls

- Augment libmpi.so, libc.so ... with libSPM.so

Second Parallel Capability

Declaration + Definition of Pclosure

```
@spm.util.dassert(predicateCb = spm.sys.sstat.amOffline)
@spm.util.dassert(predicateCb = spm.sys.pstat.amHub)
def __init():
    return spm.pclosure.macro.papply.template.openMPI.\
        policyA.defun(signature = 'signature::Hub',
                      stageCb = __taskStat,
                      );

__pc = __init();
```

Second Parallel Capability

Population + Invocation of Pclosure

```
@spm.util.dassert(predicateCb = spm.sys.sstat.amOffline)
@spm.util.dassert(predicateCb = spm.sys.pstat.amHub)
def main(pool,
         taskApiArgs,
         taskTimeout):
    # Initialize 'stage0'.
    __pc.stage0.init.main(typedef = ...);
    hdl = __pc.stage0.payload.tie();
    # Populate the template task
    hdl.spm.meta.label = '***'; # Not interested.
    hdl.spm.meta.apiArgs = taskApiArgs;
    hdl.spm.meta.timeout = taskTimeout;
    # Invoke the pmanager
    __pc.stage0.event.manage(pool = pool,
                             nSpokesMin = ...,
                             nSpokesMax = ...,
                             timeoutWaitForSpokes = ...,
                             timeoutExecution = ...
                             );
    return;
```


Second Parallel Capability

Typedef for Template Task

```
r"""
task<template>      ::struct {
  # SPM component ...
  spm               ::struct {
    meta            ::struct {
      label         ::scalar<stringSnippet> = deferred;
      apiArgs       ::dict<string,mixed>    = deferred;
      timeout       ::scalar<timeout>      = deferred;
    };

    core            ::struct {
      relaunchPre   ::scalar<bool>          = None;
      relaunchPost  ::scalar<bool>          = None;
      nameHost      ::scalar<auto>         = None;
      whoAmI        ::scalar<auto>         = None;
    };

    stat            ::struct {
      exception     ::scalar<auto>         = None;
      returnValue   ::scalar<record>       = None;
    };
  };
  # non-SPM component ...
};
"""
```

Second Parallel Capability

Callback for Status Reports

```
@spm.util.dassert (predicateCb = spm.sys.sstat.amOnline)
@spm.util.dassert (predicateCb = spm.sys.pstat.amHub)
def __taskStat (pc):
    try:
        hdl          = pc.stagel.payload.tie();
        returnValue = hdl.spm.stat.returnValue;
        if (returnValue.Has (attr = 'stdOut')):
            print ("\tstdOut   : %s", returnValue.stdOut);
        if (returnValue.Has (attr = 'stdErr')):
            print ("\tstdErr   : %s", returnValue.stdErr);
        if (returnValue.Has (attr = 'stdOutErr')):
            print ("\tstdOutErr: %s", returnValue.stdOutErr);
    except (SPMTaskDropped,
            SPMTaskLoad,
            SPMTaskEval,
            ), (hdl,):
        pass;

    return (pc.stagel.event.done(),
            None,
            )[-1];
```

SPM.Python Session

```
GNU/Linux [] spm.3.111116.trial.A.python
```



```
● >>> import pool
● >>> import demo
● >>> import os;
● >>> taskApiArgs = \
●     dict(app           = os.getcwd() + '/hello_world',
●           appOptions = "-prefix='app'",
●           );
● >>> taskTimeout = spm.util.timeout.after(seconds = 10);
✓ >>> demo.main(pool           = pool.intraAll() ,
●             taskApiArgs = taskApiArgs,
●             taskTimeout = taskTimeout)
● #: MetaStatus (hub): Waiting - ForSpokes ...
● #: MetaStatus (hub): Tasks   - Eval
●     app => 0
●     app => 1
● #: MetaStatus (hub): Tasks   - EvalDone
✓ >>> demo.main(pool           = pool.intraOnePerServer() ,
●             taskApiArgs = taskApiArgs,
●             taskTimeout = taskTimeout)
● #: MetaStatus (hub): Waiting - ForSpokes ...
● #: MetaStatus (hub): Tasks   - Eval
● #: MetaStatus (hub): Tasks   - EvalDone
● >>> exit()
GNU/Linux []
```

Conclusion

Prototyping should be frictionless



Must use original OpenMPI application

- original source code
- original binary



Original OpenMPI application must inherit support for:

- fault tolerance
- timeout
- detecting deadlocks



GNU/Linux []

~~spm.python ...~~
~~mpixn ...~~

`./hello_world -prefix "api"`

Conclusion (Cont'd)

<http://www.mbasciences.com>

{
SPM.Python distribution
Technical Briefs
Parallel Management Patterns
}

Elementary Parallel Primitives

- Clone
 - Once
 - Repeat
- Partition
 - DAG
 - List
- PartitionAggregate
 - Centralized
 - Decentralized

HPC Parallel Primitives

- Partition
 - Grid/OpenMPI

Limited Beta
Nov 30

Data / Graph Parallel Primitives

- Partition
 - Data Flow
 - Graph

Stanford U
Dec 6