

Building workflows with Spotify's



Workshop for e-Infrastructures for Massively Parallel Sequencing
Uppsala, Jan 19-20, 2015

Samuel Lampa
UU/Dept of Pharm. Biosci. & BILS
samuel.lampa@bils.se
samuel.lampa.co



Luigi Short Facts

- **Batch** (No streaming)
- Both **Command line** and **Hadoop** execution
- Does not replace Hive, Pig, Cascading etc
 - instead used to stitch many such jobs together
- Powers 1000:s of jobs every day at Spotify
- Communication / synchronization between tasks via “targets” (Normal file, HDFS, database ...)
- Dependencies hard coded in tasks
- **Pull based** (ask for task, and it figures out dependencies)
- Main features: Scheduling, Dependency Graph, Basic multi-node support (via central planner daemon)
- github.com/spotify/luigi
- Easy to install: **`pip install luigi [tornado]`**



A Basic Luigi Task

```
1 class ATask(luigi.Task):
2
3     some_param = luigi.Parameter()
4
5     def requires(self):
6         # Defines what is the upstream task to this one.
7         # Is later used by the input() function.
8         return APreviousTask()
9
10    def output(self):
11        return luigi.LocalTarget(self.input().path + ".new_extension_" + self.some_param)
12
13    def run(self):
14        # Loop over the input file, and just write each line to the output file.
15        with self.input().open() as infile, self.output().open() as outfile:
16            for line in infile:
17                outfile.writeline(line)
18
```



A Basic Luigi Task

```
1 class ATask(luigi.Task):
2
3     some_param = luigi.Parameter()
4
5     def requires(self):
6         # Defines what is the upstream task to this one.
7         # Is later used by the input() function.
8         return APreviousTask()
9
10    def output(self):
11        return luigi.LocalTarget(self.input().path + ".new_extension_" + self.some_param)
12
13    def run(self):
14        # Loop over the input file, and just write each line to the output file.
15        with self.input().open() as infile, self.output().open() as outfile:
16            for line in infile:
17                outfile.writeline(line)
18
```

The diagram illustrates the dependencies between code elements in the `ATask` class:

- An orange arrow points from `self.some_param` in the `output` method to `some_param` in the class attribute.
- A green arrow points from `self.input()` in the `requires` method to `self.input()` in the `output` method.
- A purple arrow points from `self.output()` in the `run` method to `self.output()` in the `output` method.



Calling tasks from the commandline

```
$ python mytasks.py --local-scheduler ATask --some-param SomeValue
```

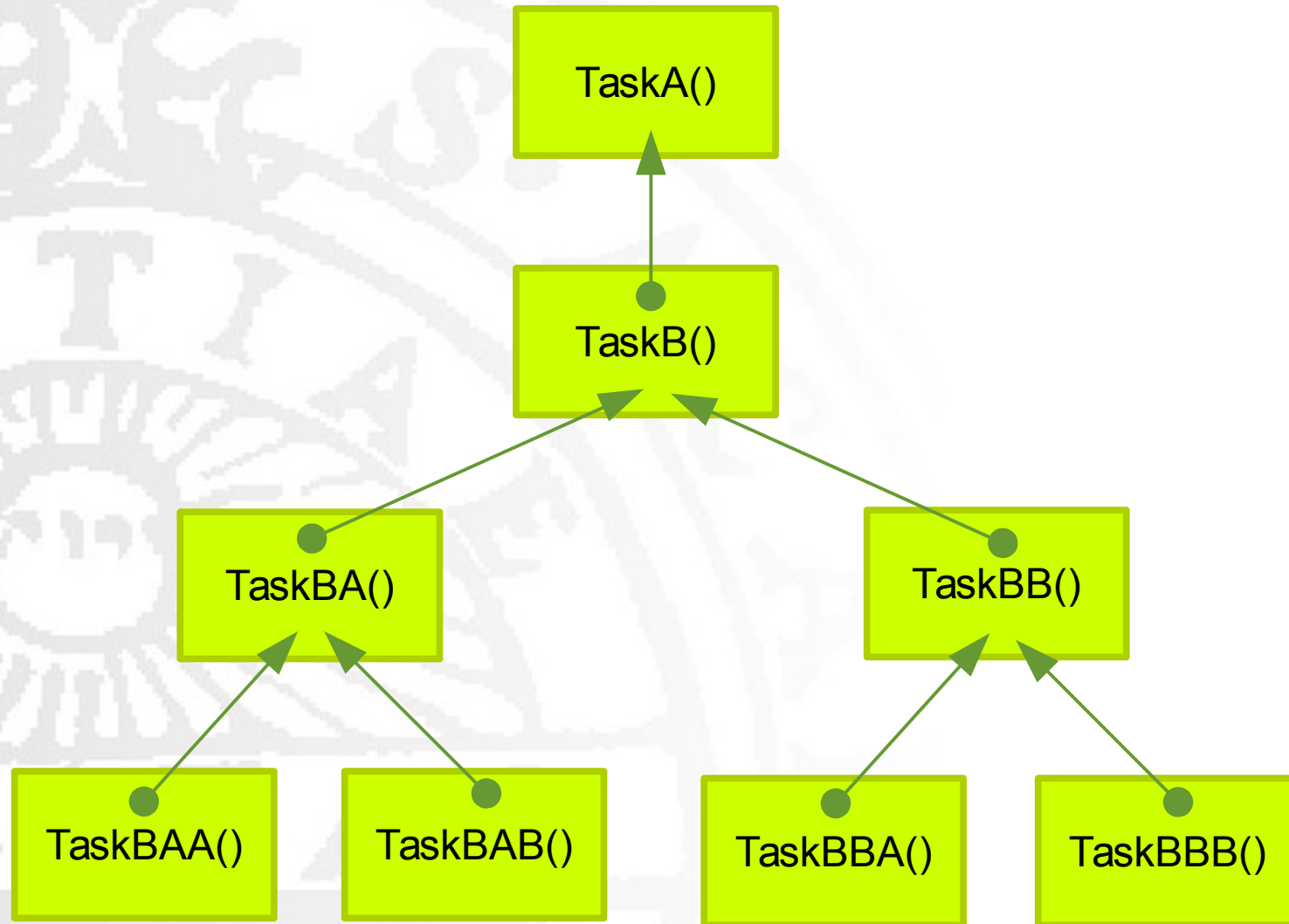


Defining workflows: Default way

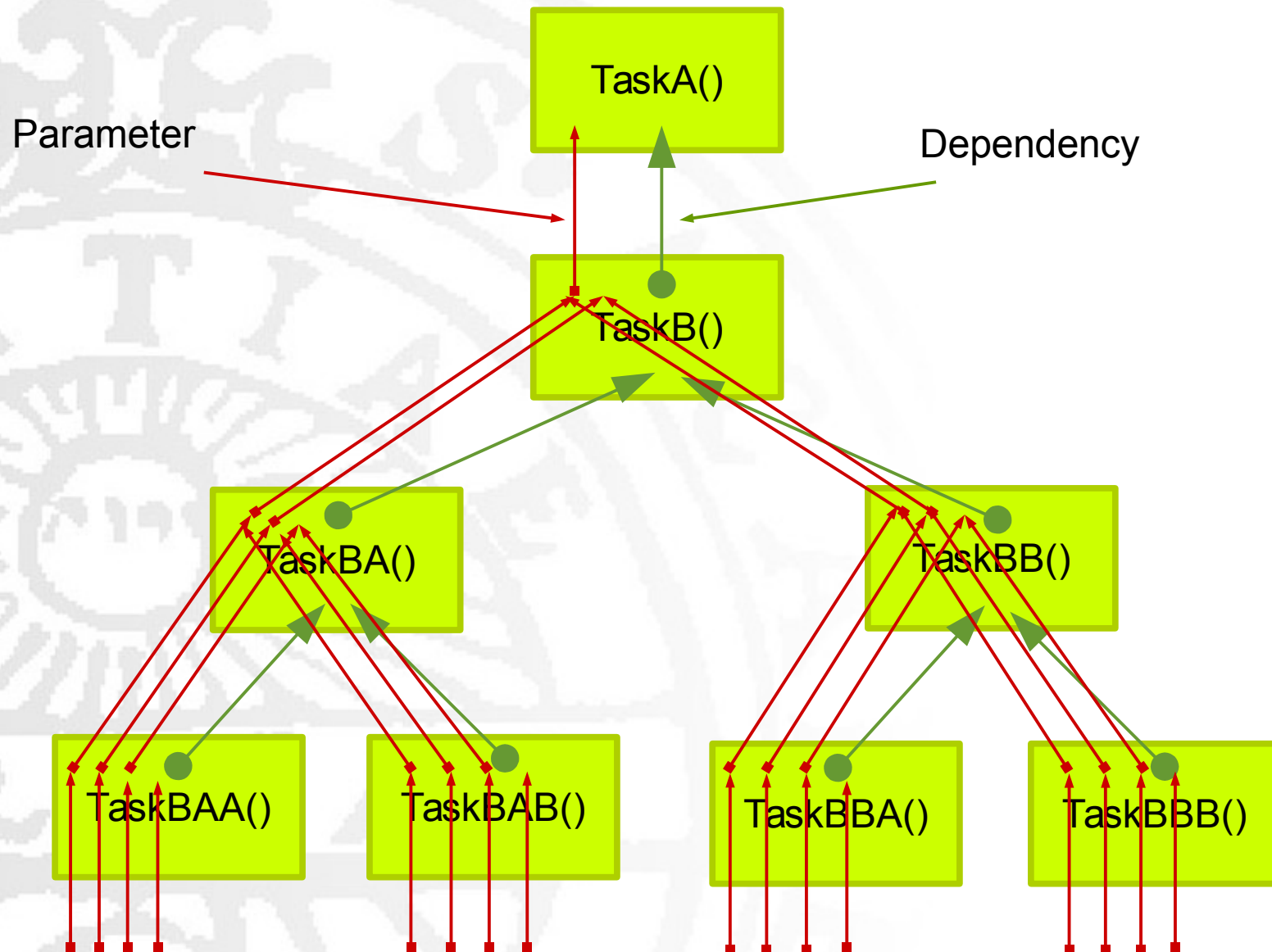
```
1 import luigi
2
3 # A luigi task
4 class HelloWorld(luigi.Task):
5     def requires(self):
6         return None
7     def output(self):
8         return luigi.LocalTarget('helloworld.txt')
9     def run(self):
10         with self.output().open('w') as outfile:
11             outfile.write('Hello World!\n')
12
13 # Another task, that depends on the above one
14 class NameSubstituter(luigi.Task):
15     name = luigi.Parameter()
16
17     def requires(self):
18         return HelloWorld() # <-- Dependency definition
19     def output(self):
20         return luigi.LocalTarget(self.input().path + '.name_' + self.name)
21     def run(self):
22         with self.input().open() as infile, self.output().open('w') as outfile:
23             text = infile.read()
24             text = text.replace('World', self.name)
25             outfile.write(text)
26
27 if __name__ == '__main__':
28     luigi.run()
29
```



Dependencies: Default way



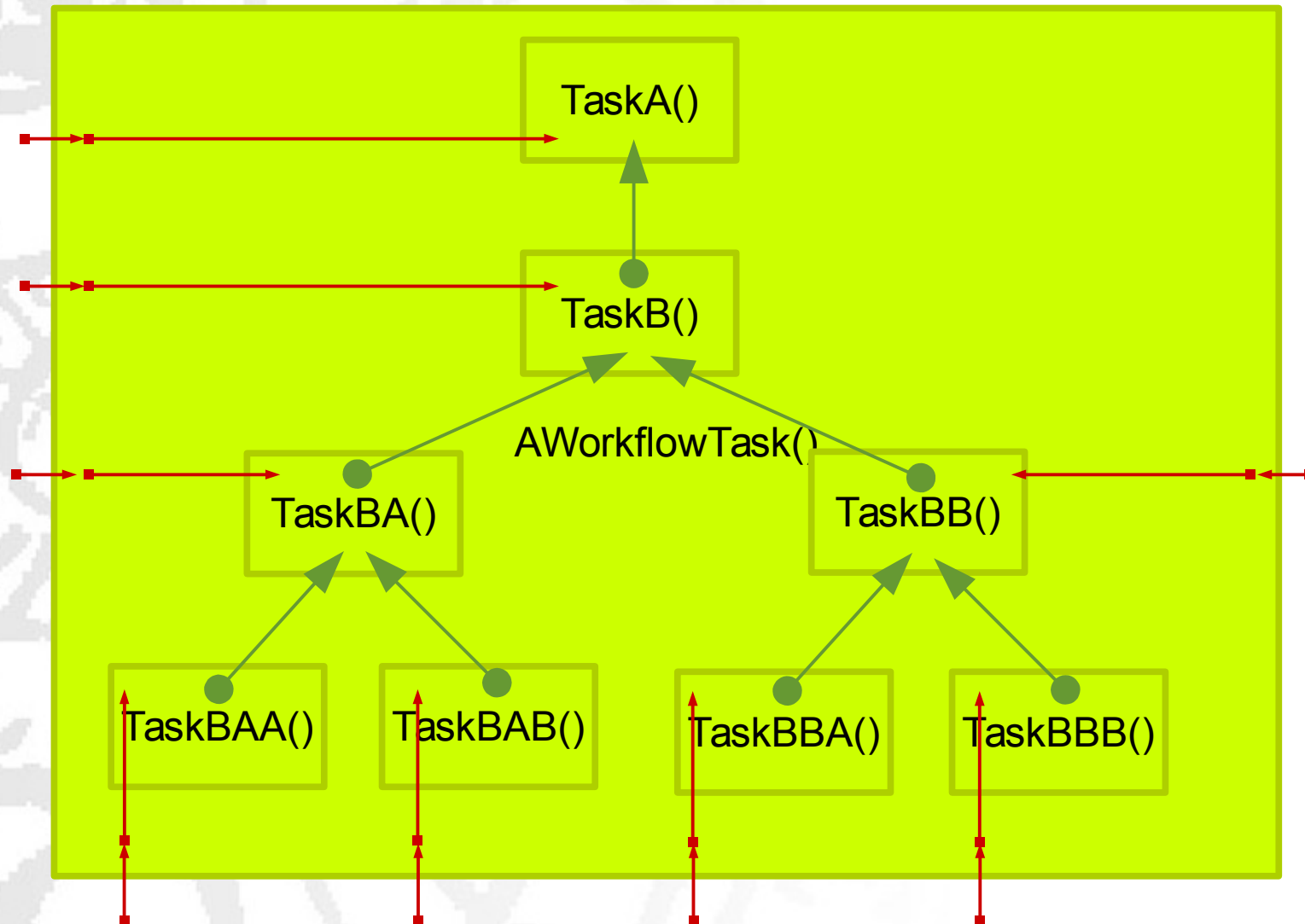
Parameters: Default way



Growing list of parameters

Growing list of parameters

Dependencies and parameters: A better way



Supporting separate network definition

```
53 class DependencyMetaTask(luigi.Task):
54     # METHODS FOR AUTOMATING DEPENDENCY MANAGEMENT
55     def get_upstream_targets(self):
56         upstream_tasks = []
57         for param_val in self.param_args:
58             if type(param_val) is dict:
59                 if 'upstream' in param_val:
60                     upstream_tasks.append(param_val['upstream']['task'])
61         return upstream_tasks
62
63     def requires(self):
64         return self.get_upstream_targets()
65
66     def get_input(self, input_name):
67         param = self.param_kwargs[input_name]
68         if type(param) is dict and 'upstream' in param:
69             return param['upstream']['task'].output()[param['upstream']['port']]
70         else:
71             return param
72
73     def get_value(self, input_name):
74         param = self.param_kwargs[input_name]
75         if type(param) is dict and 'upstream' in param:
76             input_target = param['upstream']['task'].output()[param['upstream']['port']]
77             if os.path.isfile(input_target.path):
78                 with input_target.open() as infile:
79                     csv_reader = csv.reader(infile)
80                     for row in csv_reader:
81                         if row[0] == param['upstream']['key']:
82                             return row[1]
83             else:
84                 return 'NA'
85         else:
86             return param
87
```



Using this functionality in tasks

```
1 class CreateSparseTrainDataset(DependencyMetaTask, TaskHelpers, DatasetNameMixin, AuditTrailMixin):
2
3     # INPUT TARGETS
4     train_dataset_target = luigi.Parameter()
5
6     # TASK PARAMETERS
7     replicate_id = luigi.Parameter()
8
9     # DEFINE OUTPUTS
10    def output(self):
11        basepath = self.get_input('train_dataset_target').path
12        return { "sparse_train_dataset" : luigi.LocalTarget(basepath + ".csr"),
13                "signatures" : luigi.LocalTarget(basepath + ".signatures"),
14                "log" : luigi.LocalTarget(basepath + ".csr.log") }
15
16    # WHAT THE TASK DOES
17    def run(self):
18        self.x([JAVA_PATH, "-jar jars/CreateSparseDataset.jar",
19                "-inputfile", self.get_input('train_dataset_target').path,
20                "-datasetfile", self.output()['sparse_train_dataset'].path,
21                "-signaturesoutfile", self.output()["signatures"].path,
22                "-silent"])
```



Using this new functionality in workflows, and wrapping whole workflows in Tasks

```
7 class MMSVMWorkflow(luigi.Task):
8     '''
9     This class runs the MM Workflow using Support Vector Machine
10    as the method for doing machine learning
11    '''
12
13    # WORKFLOW PARAMETERS
14    dataset_name = luigi.Parameter()
15    replicate_id = luigi.Parameter()
16    test_size = luigi.Parameter()
17    train_size = luigi.Parameter()
18    sampling_seed = luigi.Parameter(default=None)
19    sampling_method = luigi.Parameter()
20    svm_gamma = luigi.Parameter()
21    svm_cost = luigi.Parameter()
22    accounted_project = luigi.Parameter()
23    parallel_svm_train = luigi.BooleanParameter()
24    #folds_count = luigi.Parameter()
25    svm_type = luigi.Parameter()
26    svm_kernel_type = luigi.Parameter()
27
28    def __init__(self, *args, **kwargs):
29        super(MMSVMWorkflow, self).__init__(*args, **kwargs)
30        '''
31        The dependency graph is defined here!
32        '''
33
34        # -----
35        self.existing_smiles = ExistingSmiles(
36            dataset_name = self.dataset_name,
37            replicate_id = self.replicate_id,
38            accounted_project = self.accounted_project)
39
40        # -----
41        self.gen_sign_filter_subst = GenerateSignaturesFilterSubstances(smiles_target=
42            { 'upstream' : { 'task' : self.existing_smiles,
43                            'port' : 'smiles' } },
44            min_height = 1,
45            max_height = 3,
46            dataset_name = self.dataset_name,
47            replicate_id = self.replicate_id,
48            accounted_project = self.accounted_project)
```

Unit testing Luigi tasks is easy

```
2 import os
3 import time
4
5 class TestConcatenate2Files():
6     file1_path = '/tmp/luigi_concat2files_file1'
7     file2_path = '/tmp/luigi_concat2files_file2'
8
9     file1_content = 'A'*80 + '\n'
10    file2_content = 'B'*80 + '\n'
11
12    def setup(self):
13        with open(self.file1_path, 'w') as file1:
14            file1.write(self.file1_content)
15        with open(self.file2_path, 'w') as file2:
16            file2.write(self.file2_content)
17
18        self.concat2files = Concatenate2Files(
19            replicate_id='TESTID',
20            accounted_project='b2015002',
21            file1_target=luigi.LocalTarget(self.file1_path),
22            file2_target=luigi.LocalTarget(self.file2_path),
23            skip_file1_header=False,
24            skip_file2_header=False
25        )
26
27    def teardown(self):
28        os.remove(self.file1_path)
29        os.remove(self.file2_path)
30        os.remove(self.concat2files.output()['concatenated_file'].path)
31
32    def test_run(self):
33        # Run the task with a luigi worker
34        w = luigi.worker.Worker()
35        w.add(self.concat2files)
36        w.run()
37        w.stop()
38
39        with open(self.concat2files.output()['concatenated_file'].path) as concat_file:
40            concatenated_content = concat_file.read()
41
42        assert concatenated_content == self.file1_content + self.file2_content
43
```



Lessons Learned (April 2014)

- Only “pull” or “call/return” semantics, no “push” or “auto triggering by inputs”
- Not fully separate workflow
- Multiple inputs / outputs somewhat error-prone
- Dynamic typed language
- Check out for max processes limits
- No real distribution of compute or data (only common scheduling)



Lessons Learned (Update Jan 2015)

- Only “pull” or “call/return” semantics, no “push” or “auto triggering by inputs”
- Not fully separate workflow def.
 - **We found a way around this.**
- Multiple inputs / outputs somewhat error-prone
 - **We're looking into solving this too.**
- Dynamic typed language
 - **We've found that unit testing is easy.**
- Check out for max processes limits
- No real distribution of compute or data (only common scheduling)
 - **We have found ourselves running everything through SLURM anyway**





Thank you!

