

Migrating Workflows to Pegasus 5.0

Ryan Tanaka
Programmer Analyst
tanaka@isi.edu

You've installed Pegasus 5.0, next steps?

- Existing workflows written in the DAX3 format need to be updated to the new YAML format
 - Support for older formats will end in Pegasus 5.1 (~ 9 months from now)
- New workflows should be written using the new Python API
- Submit node must have Python3.5+
- Perform database upgrade

Overview of Changes Affecting Workflow Development

- File Format Changes
- default data staging configuration is **condorio**
 - If using the default prior to v4.9, **sharedfs**, this needs to be explicitly set in your properties file
`pegasus.data.configuration sharedfs`
- by default, output files have `register_replica=True`
- Default site catalog containing 2 sites is generated if none is given:
 - **local**
 - **condorpool**

Entity	Pegasus 4.x	Pegasus 5.0
Workflow	DAX (XML)	YAML
Site Catalog	XML	YAML
Transformation Catalog	Text	YAML
Replica Catalog	Text	YAML

Overview of Changes Affecting Workflow Development

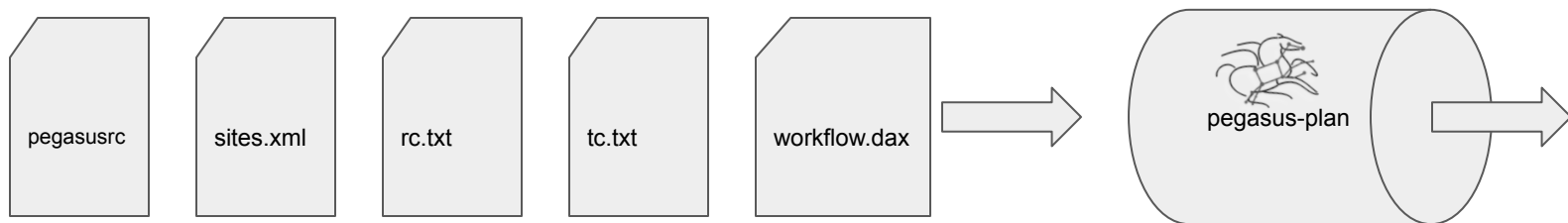
- Database upgrade required using `pegasus-db-admin update -a`
- Python API (formerly DAX3) reworked from the ground up
 - build properties file and catalogs programmatically
 - API usage simplified (adding inputs/outputs takes less lines of code)
 - plan/submit/analyze/monitor directly from workflow script

Outline

- ~~Overview of Changes~~
- **Pegasus 4.X vs Pegasus 5.0 Workflow Scripts**
- Migrating to New Python API
- Useful Tips
- Resources

Pegasus 4.x Workflow Script

- Shell script that invokes pegasus-plan at the end
- Catalogs generated either in shell script or already on filesystem
 - Transformation catalog and replica catalog may be omitted as they can be embedded into the workflow with limited functionality
- pegasusrc generated either in shell script or already on filesystem
- Workflow generation script typically invoked from within shell script to generate some <workflow name>.dax file



Pegasus 4.x Workflow Script

plan.sh

```
1 #!/bin/bash
2
3 set -e
4 set -v
5
6 TOP_DIR=$( cd "$(dirname "${BASH_SOURCE[0]}")"; pwd -P )
7 cd $TOP_DIR
8
9 # build the dax generator
10 export CLASSPATH=.:`pegasus-config --classpath`
11 javac --release 8 BlackDiamondDAX.java
12
13 # generate the dax
14 java BlackDiamondDAX /usr blackdiamond.xml
15
16 # create the site catalog
17 cat >sites.xml <<EOF
18 <?xml version="1.0" encoding="UTF-8"?>
19 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog" xmlns:xsi="http://www.w3.org/2001/
20 ecatalog http://pegasus.isi.edu/schema/sc-4.0.xsd" version="4.0">
21   <site handle="condorpool" arch="x86_64" os="LINUX" osrelease="" osversion="" glibc="">
22     <profile namespace="env" key="PEGASUS_HOME" >/usr</profile>
23     <profile namespace="condor" key="universe" >vanilla</profile>
24     <profile namespace="pegasus" key="style" >condor</profile>
25   </site>
26   <site handle="local" arch="x86_64" os="LINUX" osrelease="" osversion="7" glibc="">
27     <directory path="$HOME/workflows/scratch" type="shared-scratch" free-size="" total-size="">
28       <file-server operation="all" url="file://$HOME/workflows/scratch">
29       </file-server>
30     </directory>
31     <directory path="$HOME/workflows/outputs" type="local-storage" free-size="" total-size="">
32       <file-server operation="all" url="file://$HOME/workflows/outputs">
33       </file-server>
34     </directory>
35   </site>
36 </sitecatalog>
37 EOF
38
39 # plan the workflow
40 pegasus-plan \
41   --output-sites local \
42   --conf pegasusrc \
43   blackdiamond.xml
```

Pegasus 4.x Workflow Script

setup

plan.sh

```
1 #!/bin/bash
2
3 set -e
4 set -v
5
6 TOP_DIR=$( cd "$(dirname "${BASH_SOURCE[0]}")"; pwd -P )
7 cd $TOP_DIR
8
9 # build the dax generator
10 export CLASSPATH=.:`pegasus-config --classpath`
11 javac --release 8 BlackDiamondDAX.java
12
13 # generate the dax
14 java BlackDiamondDAX /usr blackdiamond.xml
15
16 # create the site catalog
17 cat >sites.xml <<EOF
18 <?xml version="1.0" encoding="UTF-8"?>
19 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog" xmlns:xsi="http://www.w3.org/2001/
20 ecatalog http://pegasus.isi.edu/schema/sc-4.0.xsd" version="4.0">
21   <site handle="condorpool" arch="x86_64" os="LINUX" osrelease="" osversion="" glibc="">
22     <profile namespace="env" key="PEGASUS_HOME" >/usr</profile>
23     <profile namespace="condor" key="universe" >vanilla</profile>
24     <profile namespace="pegasus" key="style" >condor</profile>
25   </site>
26   <site handle="local" arch="x86_64" os="LINUX" osrelease="" osversion="7" glibc="">
27     <directory path="$HOME/workflows/scratch" type="shared-scratch" free-size="" total-size="">
28       <file-server operation="all" url="file://$HOME/workflows/scratch">
29       </file-server>
30     </directory>
31     <directory path="$HOME/workflows/outputs" type="local-storage" free-size="" total-size="">
32       <file-server operation="all" url="file://$HOME/workflows/outputs">
33       </file-server>
34     </directory>
35   </site>
36 </sitecatalog>
37 EOF
38
39 # plan the workflow
40 pegasus-plan \
41   --output-sites local \
42   --conf pegasusrc \
43   blackdiamond.xml
```


Pegasus 4.x Workflow Script

plan.sh

setup

generate workflow file

```
1 #!/bin/bash
2
3 set -e
4 set -v
5
6 TOP_DIR=$( cd "$(dirname "${BASH_SOURCE[0]}")"; pwd -P )
7 cd $TOP_DIR
8
9 # build the dax generator
10 export CLASSPATH=.:`pegasus-config --classpath`
11 javac --release 8 BlackDiamondDAX.java
12
13 # generate the dax
14 java BlackDiamondDAX /usr blackdiamond.xml
15
16 # create the site catalog
17 cat >sites.xml <<EOF
18 <?xml version="1.0" encoding="UTF-8"?>
19 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog" xmlns:xsi="http://www.w3.org/2001/
20 ecatalog http://pegasus.isi.edu/schema/sc-4.0.xsd" version="4.0">
21 <site handle="condorpool" arch="x86_64" os="LINUX" osrelease="" osversion="" glibc="">
22 <profile namespace="env" key="PEGASUS_HOME" >/usr</profile>
23 <profile namespace="condor" key="universe" >vanilla</profile>
24 <profile namespace="pegasus" key="style" >condor</profile>
25 </site>
26 <site handle="local" arch="x86_64" os="LINUX" osrelease="" osversion="7" glibc="">
27 <directory path="$HOME/workflows/scratch" type="shared-scratch" free-size="" total-size="">
28 <file-server operation="all" url="file://$HOME/workflows/scratch">
29 </directory>
30 <directory path="$HOME/workflows/outputs" type="local-storage" free-size="" total-size="">
31 <file-server operation="all" url="file://$HOME/workflows/outputs">
32 </file-server>
33 </directory>
34 </site>
35 </sitecatalog>
36 EOF
37
38 # plan the workflow
39 pegasus-plan \
40 --output-sites local \
41 --conf pegasusr \
42 blackdiamond.xml
```

Diamond.java

```
1 /**
2  * Copyright 2007-2008 University Of Southern California
3  *
4  * Licensed under the Apache License, Version 2.0 (the "License");
5  * you may not use this file except in compliance with the License.
6  * You may obtain a copy of the License at
7  *
8  * http://www.apache.org/licenses/LICENSE-2.0
9  *
10  * Unless required by applicable law or agreed to in writing,
11  * software distributed under the License is distributed on an "AS IS" BASIS,
12  * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13  * See the License for the specific language governing permissions and
14  * limitations under the License.
15  */
16
17 import edu.isi.pegasus.planner.dax.*;
18
19 public class BlackDiamondDAX {
20
21     /**
22      * Create an example DIAMOND DAX
23      * @param args
24      */
25     public static void main(String[] args) {
26         if (args.length != 2) {
27             System.out.println("Usage: java ADAG <site_handle> <pegasus_location>");
28             System.exit(1);
29         }
30
31         try {
32             Diamond(args[0]).writeToFile(args[1]);
33         }
34         catch (Exception e) {
35             e.printStackTrace();
36         }
37     }
38
39     private static ADAG Diamond(String pegasus_location) throws Exception {
40
41         ADAG dax = new ADAG("blackdiamond");
42
43         File fb1 = new File("f.b1");
44         File fb2 = new File("f.b2");
45         File fc1 = new File("f.c1");
46         File fc2 = new File("f.c2");
47         File fd = new File("f.d");
48         File fe1 = new File("f.e1");
49         File fe2 = new File("f.e2");
50         fd.setRegister(true);
51     }
52 }
```

Pegasus 4.x Workflow Script

plan.sh

setup

generate workflow file

generate catalogs/pegasusrc

```
1 #!/bin/bash
2
3 set -e
4 set -v
5
6 TOP_DIR=$( cd "$(dirname "${BASH_SOURCE[0]}")"; pwd -P )
7 cd $TOP_DIR
8
9 # build the dax generator
10 export CLASSPATH=.:`pegasus-config --classpath`
11 javac --release 8 BlackDiamondDAX.java
12
13 # generate the dax
14 java BlackDiamondDAX /usr blackdiamond.xml
15
16 # create the site catalog
17 cat >sites.xml <<EOF
18 <?xml version="1.0" encoding="UTF-8"?>
19 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog" xmlns:xsi="http://www.w3.org/2001/
20 ecatalog http://pegasus.isi.edu/schema/sc-4.0.xsd" version="4.0">
21 <site handle="condorpool" arch="x86_64" os="LINUX" osrelease="" osversion="" glibc="">
22   <profile namespace="env" key="PEGASUS_HOME" >/usr</profile>
23   <profile namespace="condor" key="universe" >vanilla</profile>
24   <profile namespace="pegasus" key="style" >condor</profile>
25 </site>
26 <site handle="local" arch="x86_64" os="LINUX" osrelease="" osversion="7" glibc="">
27   <directory path="$HOME/workflows/scratch" type="shared-scratch" free-size="" total-size="">
28     <file-server operation="all" url="file://$HOME/workflows/scratch">
29
30     </file-server>
31   </directory>
32   <directory path="$HOME/workflows/outputs" type="local-storage" free-size="" total-size="">
33     <file-server operation="all" url="file://$HOME/workflows/outputs">
34
35     </file-server>
36   </directory>
37 </site>
38 </sitecatalog>
39 EOF
40
41 # plan the workflow
42 pegasus-plan \
43   --output-sites local \
44   --conf pegasusrc \
45   blackdiamond.xml
```

Diamond.java

```
1 /**
2  * Copyright 2007-2008 University Of Southern California
3  *
4  * Licensed under the Apache License, Version 2.0 (the "License");
5  * you may not use this file except in compliance with the License.
6  * You may obtain a copy of the License at
7  *
8  * http://www.apache.org/licenses/LICENSE-2.0
9  *
10  * Unless required by applicable law or agreed to in writing,
11  * software distributed under the License is distributed on an "AS IS" BASIS,
12  * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13  * See the License for the specific language governing permissions and
14  * limitations under the License.
15  */
16
17 import edu.isi.pegasus.planner.dax.*;
18
19 public class BlackDiamondDAX {
20
21   /**
22    * Create an example DIAMOND DAX
23    * @param args
24    */
25   public static void main(String[] args) {
26     if (args.length != 2) {
27       System.out.println("Usage: java ADAG <site_handle> <pegasus_location>");
28       System.exit(1);
29     }
30
31     try {
32       Diamond(args[0]).writeToFile(args[1]);
33     }
34     catch (Exception e) {
35       e.printStackTrace();
36     }
37   }
38
39   private static ADAG Diamond(String pegasus_location) throws Exception {
40
41     ADAG dax = new ADAG("blackdiamond");
42
43     File fb1 = new File("f.b1");
44     File fb2 = new File("f.b2");
45     File fc1 = new File("f.c1");
46     File fc2 = new File("f.c2");
47     File fd = new File("f.d");
48
49     File fe1 = new File("f.e1");
50     File fe2 = new File("f.e2");
51     fd.setRegister(true);
```

Pegasus 4.x Workflow Script

plan.sh

setup

generate workflow file

generate catalogs/pegasusrc

invoke pegasus-plan

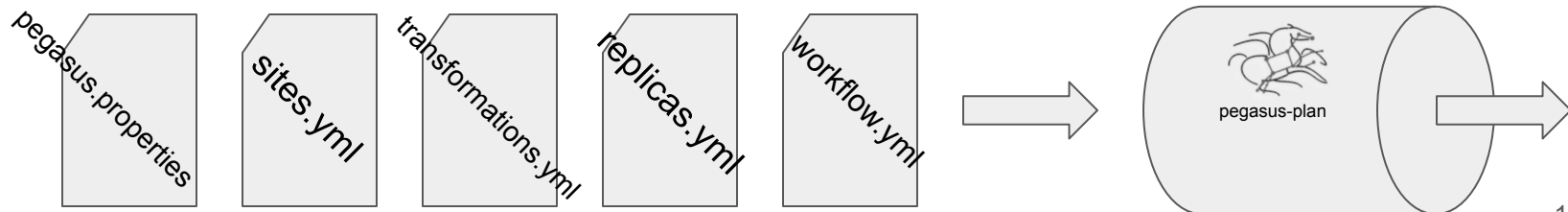
```
1 #!/bin/bash
2
3 set -e
4 set -v
5
6 TOP_DIR=$( cd "$(dirname "${BASH_SOURCE[0]}")"; pwd -P )
7 cd $TOP_DIR
8
9 # build the dax generator
10 export CLASSPATH=.:`pegasus-config --classpath`
11 javac --release 8 BlackDiamondDAX.java
12
13 # generate the dax
14 java BlackDiamondDAX /usr blackdiamond.xml
15
16 # create the site catalog
17 cat >sites.xml <<EOF
18 <?xml version="1.0" encoding="UTF-8"?>
19 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog" xmlns:xsi="http://www.w3.org/2001/
20 ecatalog http://pegasus.isi.edu/schema/sc-4.0.xsd" version="4.0">
21 <site handle="condorpool" arch="x86_64" os="LINUX" osrelease="" osversion="" glibc="">
22   <profile namespace="env" key="PEGASUS_HOME" >/usr</profile>
23   <profile namespace="condor" key="universe" >vanilla</profile>
24   <profile namespace="pegasus" key="style" >condor</profile>
25 </site>
26 <site handle="local" arch="x86_64" os="LINUX" osrelease="" osversion="7" glibc="">
27   <directory path="$HOME/workflows/scratch" type="shared-scratch" free-size="" total-size="">
28     <file-server operation="all" url="file://$HOME/workflows/scratch">
29       </file-server>
30     </directory>
31     <directory path="$HOME/workflows/outputs" type="local-storage" free-size="" total-size="">
32       <file-server operation="all" url="file://$HOME/workflows/outputs">
33         </file-server>
34       </directory>
35     </site>
36 </sitecatalog>
37 EOF
38
39 # plan the workflow
40 pegasus-plan \
41   --output-sites local \
42   --conf pegasusrc \
43   blackdiamond.xml
```

Diamond.java

```
1 /**
2  * Copyright 2007-2008 University Of Southern California
3  *
4  * Licensed under the Apache License, Version 2.0 (the "License");
5  * you may not use this file except in compliance with the License.
6  * You may obtain a copy of the License at
7  *
8  * http://www.apache.org/licenses/LICENSE-2.0
9  *
10  * Unless required by applicable law or agreed to in writing,
11  * software distributed under the License is distributed on an "AS IS" BASIS,
12  * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13  * See the License for the specific language governing permissions and
14  * limitations under the License.
15  */
16
17 import edu.isi.pegasus.planner.dax.*;
18
19 public class BlackDiamondDAX {
20
21   /**
22    * Create an example DIAMOND DAX
23    * @param args
24    */
25   public static void main(String[] args) {
26     if (args.length != 2) {
27       System.out.println("Usage: java ADAG <site_handle> <pegasus_location>");
28       System.exit(1);
29     }
30
31     try {
32       Diamond(args[0]).writeToFile(args[1]);
33     }
34     catch (Exception e) {
35       e.printStackTrace();
36     }
37   }
38
39   private static ADAG Diamond(String pegasus_location) throws Exception {
40
41     ADAG dax = new ADAG("blackdiamond");
42
43     File fb1 = new File("f.b1");
44     File fb2 = new File("f.b2");
45     File fc1 = new File("f.c1");
46     File fc2 = new File("f.c2");
47     File fd = new File("f.d");
48
49     File fe1 = new File("f.e1");
50     File fe2 = new File("f.e2");
51     fd.setRegister(true);
```

Pegasus 5.0 Workflow Script

- Single script using the Pegasus 5.0 Python API
- Programmatically generate:
 - pegasus.properties (you may have seen pegasusrc, same file)
 - Site Catalog (sites.yml)
 - default one will be used if none is picked up by pegasus-plan
 - Replica Catalog (replicas.yml)
 - Transformation Catalog (transformations.yml)
- Catalogs remain separate entities from the abstract workflow
- Plan and run directly from this script



Pegasus 5.0 Workflow Script

workflow.py

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Pegasus 5.0 Workflow Script

workflow.py

generate pegasus.properties

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```


Pegasus 5.0 Workflow Script

workflow.py

generate pegasus.properties
generate Replica Catalog

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Pegasus 5.0 Workflow Script

workflow.py

generate pegasus.properties
generate Replica Catalog
generate Transformation Catalog

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Pegasus 5.0 Workflow Script

workflow.py

generate pegasus.properties
generate Replica Catalog
generate Transformation Catalog
generate workflow

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Pegasus 5.0 Workflow Script

workflow.py

generate pegasus.properties
generate Replica Catalog
generate Transformation Catalog
generate workflow
plan and run workflow

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Outline

- ~~Overview of Changes~~
- ~~Pegasus 4.X vs Pegasus 5.0 Workflow Scripts~~
- **Migrating to New Python API**
- Useful Tips
- Resources

Migrating Workflow Scripts to New Python API

- Migrate catalogs from 4.x format to YAML
 - Generate new catalogs using the Python API (recommended) **OR**
 - Convert existing catalogs using provided catalog converters
 - `pegasus-sc-converter -i sites.xml -o sites.yml`
 - `pegasus-tc-converter -i tc.txt -I Text -O YAML -o transformations.yml`
 - `pegasus-rc-converter -I File -O YAML -i rc.txt -o replicas.yml transformations.yml`
- Migrate workflow generation script from using `Pegasus.DAX3` to `Pegasus.api`

Python API: Site Catalog

sites.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog"
3   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4   xsi:schemaLocation="http://pegasus.isi.edu/schema/sitecatalog http://pegasus.isi.edu/schem
5   version="4.0">
6
7   <site handle="local" arch="x86_64" os="LINUX">
8     <directory type="shared-scratch" path="/tmp/workflows/scratch">
9       <file-server operation="all" url="file:///tmp/workflows/scratch"/>
10     </directory>
11     <directory type="local-storage" path="/tmp/workflows/outputs">
12       <file-server operation="all" url="file:///tmp/workflows/outputs"/>
13     </directory>
14   </site>
15
16   <site handle="condor_pool" arch="x86_64" os="LINUX">
17     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="auxiliary"/>
18     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="compute"/>
19     <directory type="shared-scratch" path="/lustre">
20       <file-server operation="all" url="gsiftp://smarty.isi.edu/lustre"/>
21     </directory>
22   </site>
23
24   <site handle="staging_site" arch="x86_64" os="LINUX">
25     <directory type="shared-scratch" path="/data">
26       <file-server operation="put" url="scp://obelix.isi.edu/data"/>
27       <file-server operation="get" url="http://obelix.isi.edu/data"/>
28     </directory>
29   </site>
30 </sitecatalog>
```

generate_sc.py

```
1 # create a SiteCatalog object
2 sc = SiteCatalog()
3
4 # create a "local" site
5 local = Site("local", arch=Arch.X86_64, os_type=OS.LINUX)
6
7 # create and add a shared scratch and local storage directories to the site "local"
8 local_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/tmp/workflows/scratch")\
9   .add_file_servers(FileServer("file:///tmp/workflows/scratch", Operation.ALL))
10
11 local_local_storage_dir = Directory(Directory.LOCAL_STORAGE, path="/tmp/workflows/outputs")\
12   .add_file_servers(FileServer("file:///tmp/workflows/outputs", Operation.ALL))
13
14 local.add_directories(local_shared_scratch_dir, local_local_storage_dir)
15
16 # create a "condorpool" site
17 condorpool = Site("condorpool", arch=Arch.X86_64, os_type=OS.LINUX)
18
19 # create and add job managers to the site "condorpool"
20 condorpool.add_grids(
21   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.AUXILIARY),
22   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.COMPUTE)
23 )
24
25 # create and add a shared scratch directory to the site "condorpool"
26 condorpool_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/lustre")\
27   .add_file_servers(FileServer("gsiftp://smarty.isi.edu/lustre", Operation.ALL))
28 condorpool.add_directories(condorpool_shared_scratch_dir)
29
30 # create a "staging_site" site
31 staging_site = Site("staging_site", arch=Arch.X86_64, os_type=OS.LINUX)
32
33 # create and add a shared scratch directory to the site "staging_site"
34 staging_site_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/data")\
35   .add_file_servers(
36     FileServer("scp://obelix.isi.edu/data", Operation.PUT),
37     FileServer("http://obelix.isi.edu/data", Operation.GET)
38   )
39 staging_site.add_directories(staging_site_shared_scratch_dir)
40
41 # add all the sites to the site catalog object
42 sc.add_sites(
43   local,
44   condorpool,
45   staging_site
46 )
47
48 # write the site catalog to the default path "./sites.xml"
49 sc.write()
```

Python API: Site Catalog

sites.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog"
3   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4   xsi:schemaLocation="http://pegasus.isi.edu/schema/sitecatalog http://pegasus.isi.edu/schem
5   version="4.0">
6
7   <site handle="local" arch="x86_64" os="LINUX">
8     <directory type="shared-scratch" path="/tmp/workflows/scratch">
9       <file-server operation="all" url="file:///tmp/workflows/scratch"/>
10     </directory>
11     <directory type="local-storage" path="/tmp/workflows/outputs">
12       <file-server operation="all" url="file:///tmp/workflows/outputs"/>
13     </directory>
14   </site>
15
16   <site handle="condor_pool" arch="x86_64" os="LINUX">
17     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="auxiliary"/>
18     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="compute"/>
19     <directory type="shared-scratch" path="/lustre">
20       <file-server operation="all" url="gsiftp://smarty.isi.edu/lustre"/>
21     </directory>
22   </site>
23
24   <site handle="staging_site" arch="x86_64" os="LINUX">
25     <directory type="shared-scratch" path="/data">
26       <file-server operation="put" url="scp://obelix.isi.edu/data"/>
27       <file-server operation="get" url="http://obelix.isi.edu/data"/>
28     </directory>
29   </site>
30 </sitecatalog>
```

generate_sc.py

```
1 # create a SiteCatalog object
2 sc = SiteCatalog()
3
4 # create a "local" site
5 local = Site("local", arch=Arch.X86_64, os_type=OS.LINUX)
6
7 # create and add a shared scratch and local storage directories to the site "local"
8 local_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/tmp/workflows/scratch")\
9   .add_file_servers(FileServer("file:///tmp/workflows/scratch", Operation.ALL))
10
11 local_local_storage_dir = Directory(Directory.LOCAL_STORAGE, path="/tmp/workflows/outputs")\
12   .add_file_servers(FileServer("file:///tmp/workflows/outputs", Operation.ALL))
13
14 local.add_directories(local_shared_scratch_dir, local_local_storage_dir)
15
16 # create a "condorpool" site
17 condorpool = Site("condorpool", arch=Arch.X86_64, os_type=OS.LINUX)
18
19 # create and add job managers to the site "condorpool"
20 condorpool.add_grids(
21   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.AUXILIARY),
22   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.COMPUTE)
23 )
24
25 # create and add a shared scratch directory to the site "condorpool"
26 condorpool_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/lustre")\
27   .add_file_servers(FileServer("gsiftp://smarty.isi.edu/lustre", Operation.ALL))
28 condorpool.add_directories(condorpool_shared_scratch_dir)
29
30 # create a "staging_site" site
31 staging_site = Site("staging_site", arch=Arch.X86_64, os_type=OS.LINUX)
32
33 # create and add a shared scratch directory to the site "staging_site"
34 staging_site_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/data")\
35   .add_file_servers(
36     FileServer("scp://obelix.isi.edu/data", Operation.PUT),
37     FileServer("http://obelix.isi.edu/data", Operation.GET)
38   )
39 staging_site.add_directories(staging_site_shared_scratch_dir)
40
41 # add all the sites to the site catalog object
42 sc.add_sites(
43   local,
44   condorpool,
45   staging_site
46 )
47
48 # write the site catalog to the default path "./sites.xml"
49 sc.write()
```

Python API: Site Catalog

sites.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog"
3   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4   xsi:schemaLocation="http://pegasus.isi.edu/schema/sitecatalog http://pegasus.isi.edu/schema/sitecatalog.xsd"
5   version="4.0">
6
7   <site handle="local" arch="x86_64" os="LINUX">
8     <directory type="shared-scratch" path="/tmp/workflows/scratch">
9       <file-server operation="all" url="file:///tmp/workflows/scratch"/>
10     </directory>
11     <directory type="local-storage" path="/tmp/workflows/outputs">
12       <file-server operation="all" url="file:///tmp/workflows/outputs"/>
13     </directory>
14   </site>
15
16   <site handle="condor_pool" arch="x86_64" os="LINUX">
17     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="auxiliary"/>
18     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="compute"/>
19     <directory type="shared-scratch" path="/lustre">
20       <file-server operation="all" url="gsiftp://smarty.isi.edu/lustre"/>
21     </directory>
22   </site>
23
24   <site handle="staging_site" arch="x86_64" os="LINUX">
25     <directory type="shared-scratch" path="/data">
26       <file-server operation="put" url="scp://obelix.isi.edu/data"/>
27       <file-server operation="get" url="http://obelix.isi.edu/data"/>
28     </directory>
29   </site>
30 </sitecatalog>
```

generate_sc.py

```
1 # create a SiteCatalog object
2 sc = SiteCatalog()
3
4 # create a "local" site
5 local = Site("local", arch=Arch.X86_64, os_type=OS.LINUX)
6
7 # create and add a shared scratch and local storage directories to the site "local"
8 local_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/tmp/workflows/scratch")\
9   .add_file_servers(FileServer("file:///tmp/workflows/scratch", Operation.ALL))
10
11 local_local_storage_dir = Directory(Directory.LOCAL_STORAGE, path="/tmp/workflows/outputs")\
12   .add_file_servers(FileServer("file:///tmp/workflows/outputs", Operation.ALL))
13
14 local.add_directories(local_shared_scratch_dir, local_local_storage_dir)
15
16 # create a "condorpool" site
17 condorpool = Site("condorpool", arch=Arch.X86_64, os_type=OS.LINUX)
18
19 # create and add job managers to the site "condorpool"
20 condorpool.add_grids(
21   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.AUXILIARY),
22   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.COMPUTE)
23 )
24
25 # create and add a shared scratch directory to the site "condorpool"
26 condorpool_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/lustre")\
27   .add_file_servers(FileServer("gsiftp://smarty.isi.edu/lustre", Operation.ALL))
28 condorpool.add_directories(condorpool_shared_scratch_dir)
29
30 # create a "staging_site" site
31 staging_site = Site("staging_site", arch=Arch.X86_64, os_type=OS.LINUX)
32
33 # create and add a shared scratch directory to the site "staging_site"
34 staging_site_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/data")\
35   .add_file_servers(
36     FileServer("scp://obelix.isi.edu/data", Operation.PUT),
37     FileServer("http://obelix.isi.edu/data", Operation.GET)
38   )
39 staging_site.add_directories(staging_site_shared_scratch_dir)
40
41 # add all the sites to the site catalog object
42 sc.add_sites(
43   local,
44   condorpool,
45   staging_site
46 )
47
48 # write the site catalog to the default path "./sites.xml"
49 sc.write()
```

Python API: Site Catalog

sites.xml

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <sitecatalog xmlns="http://pegasus.isi.edu/schema/sitecatalog"
3   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4   xsi:schemaLocation="http://pegasus.isi.edu/schema/sitecatalog http://pegasus.isi.edu/schema/sitecatalog.xsd"
5   version="4.0">
6
7   <site handle="local" arch="x86_64" os="LINUX">
8     <directory type="shared-scratch" path="/tmp/workflows/scratch">
9       <file-server operation="all" url="file:///tmp/workflows/scratch"/>
10     </directory>
11     <directory type="local-storage" path="/tmp/workflows/outputs">
12       <file-server operation="all" url="file:///tmp/workflows/outputs"/>
13     </directory>
14   </site>
15
16   <site handle="condor_pool" arch="x86_64" os="LINUX">
17     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="auxiliary"/>
18     <grid type="gt5" contact="smarty.isi.edu/jobmanager-pbs" scheduler="PBS" jobtype="compute"/>
19     <directory type="shared-scratch" path="/lustre">
20       <file-server operation="all" url="gsiftp://smarty.isi.edu/lustre"/>
21     </directory>
22   </site>
23
24   <site handle="staging_site" arch="x86_64" os="LINUX">
25     <directory type="shared-scratch" path="/data">
26       <file-server operation="put" url="scp://obelix.isi.edu/data"/>
27       <file-server operation="get" url="http://obelix.isi.edu/data"/>
28     </directory>
29   </site>
30 </sitecatalog>
```

generate_sc.py

```
1 # create a SiteCatalog object
2 sc = SiteCatalog()
3
4 # create a "local" site
5 local = Site("local", arch=Arch.X86_64, os_type=OS.LINUX)
6
7 # create and add a shared scratch and local storage directories to the site "local"
8 local_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/tmp/workflows/scratch")\
9   .add_file_servers(FileServer("file:///tmp/workflows/scratch", Operation.ALL))
10
11 local_local_storage_dir = Directory(Directory.LOCAL_STORAGE, path="/tmp/workflows/outputs")\
12   .add_file_servers(FileServer("file:///tmp/workflows/outputs", Operation.ALL))
13
14 local.add_directories(local_shared_scratch_dir, local_local_storage_dir)
15
16 # create a "condorpool" site
17 condorpool = Site("condorpool", arch=Arch.X86_64, os_type=OS.LINUX)
18
19 # create and add job managers to the site "condorpool"
20 condorpool.add_grids(
21   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.AUXILIARY),
22   Grid(Grid.GT5, contact="smarty.isi.edu/jobmanager-pbs", scheduler_type=Scheduler.PBS, job_type=SupportedJobs.COMPUTE)
23 )
24
25 # create and add a shared scratch directory to the site "condorpool"
26 condorpool_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/lustre")\
27   .add_file_servers(FileServer("gsiftp://smarty.isi.edu/lustre", Operation.ALL))
28 condorpool.add_directories(condorpool_shared_scratch_dir)
29
30 # create a "staging_site" site
31 staging_site = Site("staging_site", arch=Arch.X86_64, os_type=OS.LINUX)
32
33 # create and add a shared scratch directory to the site "staging_site"
34 staging_site_shared_scratch_dir = Directory(Directory.SHARED_SCRATCH, path="/data")\
35   .add_file_servers(
36     FileServer("scp://obelix.isi.edu/data", Operation.PUT),
37     FileServer("http://obelix.isi.edu/data", Operation.GET)
38   )
39 staging_site.add_directories(staging_site_shared_scratch_dir)
40
41 # add all the sites to the site catalog object
42 sc.add_sites(
43   local,
44   condorpool,
45   staging_site
46 )
47
48 # write the site catalog to the default path "./sites.xml"
49 sc.write()
```


Python API: Transformation Catalog

generate_tc.py

tc.txt

```
1 cont centos-pegasus{
2   type "docker"
3   image "docker:///rynge/montage:latest"
4   mount "/Volumes/Work/lfs1:/shared-data/:ro"
5   profile env "JAVA_HOME" "/opt/java/1.6"
6 }
7
8 tr example::keg:1.0 {
9
10  profile env "APP_HOME" "/tmp/myscratch"
11  profile env "JAVA_HOME" "/opt/java/1.6"
12
13  site isi {
14    pfn "/path/to/keg"
15    arch "x86"
16    os "linux"
17    type "INSTALLED"
18    container "centos-pegasus"
19  }
20 }
```

```
1 # create the TransformationCatalog object
2 tc = TransformationCatalog()
3
4 # create and add the centos-pegasus container
5 centos_cont = Container(
6     "centos-pegasus",
7     Container.DOCKER,
8     "docker:///rynge/montage:latest",
9     mounts=["/Volumes/Workf/lfs1:/shared-data/:ro"]
10    ).add_profiles(Namespace.ENV, JAVA_HOME="/opt/java/1.6")
11
12 tc.add_containers(centos_cont)
13
14 # create and add the transformation
15 keg = Transformation(
16     "keg",
17     namespace="example",
18     version="1.0",
19     site="isi",
20     pfn="/path/to/keg",
21     is_stageable=False,
22     container=centos_cont
23    ).add_profiles(Namespace.ENV, APP_HOME="/tmp/myscratch", JAVA_HOME="/opt/java/1.6")
24
25 tc.add_transformations(keg)
26
27 # write the transformation catalog to the default file path "./transformations.yml"
28 tc.write()
```

Python API: Transformation Catalog

generate_tc.py

tc.txt

```
1 cont centos-pegasus{
2   type "docker"
3   image "docker:///rynge/montage:latest"
4   mount "/Volumes/Work/lfs1:/shared-data/:ro"
5   profile env "JAVA_HOME" "/opt/java/1.6"
6 }
7
8 tr example::keg:1.0 {
9
10  profile env "APP_HOME" "/tmp/myscratch"
11  profile env "JAVA_HOME" "/opt/java/1.6"
12
13  site isi {
14    pfn "/path/to/keg"
15    arch "x86"
16    os "linux"
17    type "INSTALLED"
18    container "centos-pegasus"
19  }
20 }
```

```
1 # create the TransformationCatalog object
2 tc = TransformationCatalog()
3
4 # create and add the centos-pegasus container
5 centos_cont = Container(
6     "centos-pegasus",
7     Container.DOCKER,
8     "docker:///rynge/montage:latest",
9     mounts=["/Volumes/Workf/lfs1:/shared-data/:ro"]
10    ).add_profiles(Namespace.ENV, JAVA_HOME="/opt/java/1.6")
11
12 tc.add_containers(centos_cont)
13
14 # create and add the transformation
15 keg = Transformation(
16     "keg",
17     namespace="example",
18     version="1.0",
19     site="isi",
20     pfn="/path/to/keg",
21     is_stageable=False,
22     container=centos_cont
23    ).add_profiles(Namespace.ENV, APP_HOME="/tmp/myscratch", JAVA_HOME="/opt/java/1.6")
24
25 tc.add_transformations(keg)
26
27 # write the transformation catalog to the default file path "./transformations.yml"
28 tc.write()
```


Python API: Transformation Catalog

generate_tc.py

tc.txt

```
1 cont centos-pegasus{
2   type "docker"
3   image "docker:///rynge/montage:latest"
4   mount "/Volumes/Work/lfs1:/shared-data/:ro"
5   profile env "JAVA_HOME" "/opt/java/1.6"
6 }
7
8 tr example::keg:1.0 {
9
10  profile env "APP_HOME" "/tmp/myscratch"
11  profile env "JAVA_HOME" "/opt/java/1.6"
12
13  site isi {
14    pfn "/path/to/keg"
15    arch "x86"
16    os "linux"
17    type "INSTALLED"
18    container "centos-pegasus"
19  }
20 }
```

```
1 # create the TransformationCatalog object
2 tc = TransformationCatalog()
3
4 # create and add the centos-pegasus container
5 centos_cont = Container(
6     "centos-pegasus",
7     Container.DOCKER,
8     "docker:///rynge/montage:latest",
9     mounts=["/Volumes/Workf/lfs1:/shared-data/:ro"]
10    ).add_profiles(Namespace.ENV, JAVA_HOME="/opt/java/1.6")
11
12 tc.add_containers(centos_cont)
13
14 # create and add the transformation
15 keg = Transformation(
16     "keg",
17     namespace="example",
18     version="1.0",
19     site="isi",
20     pfn="/path/to/keg",
21     is_stageable=False,
22     container=centos_cont
23    ).add_profiles(Namespace.ENV, APP_HOME="/tmp/myscratch", JAVA_HOME="/opt/java/1.6")
24
25 tc.add_transformations(keg)
26
27 # write the transformation catalog to the default file path "./transformations.yml"
28 tc.write()
```

Python API: Replica Catalog

rc.txt

```
1 f.a file:///Volumes/data/inputs/f.a site="local"  
2  
3 f.b file:///Volumes/data/inputs/f.b site="local"
```

generate_rc.py

```
1 rc = ReplicaCatalog()  
2  
3 rc.add_replica(  
4     site="local",  
5     lfn="f.a",  
6     pfn="/Volumes/data/inputs/f.a"  
7 )  
8  
9 rc.add_replica(  
10     site="local",  
11     lfn="f.b",  
12     pfn="/Volumes/data/inputs/f.b"  
13 )  
14  
15 rc.write()
```

Python API: Replica Catalog

rc.txt

```
1 f.a file:///Volumes/data/inputs/f.a site="local"  
2  
3 f.b file:///Volumes/data/inputs/f.b site="local"
```

generate_rc.py

```
1 rc = ReplicaCatalog()  
2  
3 rc.add_replica(  
4     site="local",  
5     lfn="f.a",  
6     pfn="/Volumes/data/inputs/f.a"  
7 )  
8  
9 rc.add_replica(  
10     site="local",  
11     lfn="f.b",  
12     pfn="/Volumes/data/inputs/f.b"  
13 )  
14  
15 rc.write()
```

Python API: Replica Catalog

rc.txt

```
1 f.a file:///Volumes/data/inputs/f.a site="local"  
2  
3 f.b file:///Volumes/data/inputs/f.b site="local"
```

generate_rc.py

```
1 rc = ReplicaCatalog()  
2  
3 rc.add_replica(  
4     site="local",  
5     lfn="f.a",  
6     pfn="/Volumes/data/inputs/f.a"  
7 )  
8  
9 rc.add_replica(  
10     site="local",  
11     lfn="f.b",  
12     pfn="/Volumes/data/inputs/f.b"  
13 )  
14  
15 rc.write()
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```


Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```

Python API: Workflow

daxgen.py

```
1 diamond = ADAG("diamond")
2
3 a = File("f.a")
4 b1 = File("f.b1")
5 b2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     name="preprocess",
10    version="4.0"
11)
12
13 preprocess.addArguments("-T60", "-i", a, "-o", b1, b2)
14
15 preprocess.uses(a, link=Link.INPUT)
16
17 preprocess.uses(b1, link=Link.OUTPUT)
18 preprocess.uses(b2, link=Link.OUTPUT)
19
20 diamond.addJob(preprocess)
```

workflow.py

```
1 wf = Workflow("diamond")
2
3 fa = File("f.a")
4 fb1 = File("f.b1")
5 fb2 = File("f.b2")
6
7 preprocess = Job(
8     namespace="diamond",
9     transformation="preprocess",
10    version="4.0"
11)
12
13 preprocess.add_args("-T", "3", "-i", fa, "-o", fb1, fb2)
14
15 preprocess.add_inputs(fa)
16
17 preprocess.add_outputs(fb1, fb2, register_replica=True, stage_out=True)
18
19 wf.add_jobs(preprocess)
```


Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```


Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```

Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```

Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```

Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```

Python API: Hierarchical Workflow

daxgen.py

```
1 # Create a abstract dag
2 adag = ADAG('local-hierarchy')
3
4 daxfile = File('blackdiamond.dax')
5 dax1 = DAX(daxfile)
6 # DAX jobs are called with same arguments passed,
7 # while planning the root level dax
8 dax1.addArguments('--output-site local')
9 dax1.addArguments('-vvv')
10 adag.addJob(dax1)
11
12
13 # this dax job uses a pre-existing dax file
14 # that has to be present in the replica catalog
15 daxfile2 = File('sleep.dax')
16 dax2 = DAX(daxfile2)
17 dax2.addArguments('--output-site local')
18 dax2.addArguments('-vvv')
19 adag.addJob(dax2)
```

workflow.py

```
1 blackdiamond_wf = SubWorkflow("blackdiamond.yml", is_planned=False)
2 blackdiamond_wf.add_args(
3     "--input-dir",
4     "input",
5     "--output-sites",
6     "local",
7     "-vvv"
8 )
9
10 sleep_wf = SubWorkflow("sleep.yml", is_planned=False)
11 sleep_wf.add_args("--output-sites", "local", "-vvv")
12
13 wf.add_jobs(blackdiamond_wf, sleep_wf)
```

Python API: Adding Profiles

Profiles can be applied to:

- FileServer
- Site
- Container
- Transformation Site
- Transformation
- Job
- SubWorkflow
- Workflow

daxgen.py

```
1 # add a Profile Object
2 job.addProfile(Profile(Namespace.ENV, 'PATH', '/bin'))
3
4 # add a profile directly
5 job.profile(Namespace.CONDOR, "universe", "vanilla")
```


Python API: Adding Profiles

daxgen.py

```
1 # add a Profile Object
2 job.addProfile(Profile(Namespace.ENV, 'PATH', '/bin'))
3
4 # add a profile directly
5 job.profile(Namespace.CONDOR, "universe", "vanilla")
```

Create a Profile object, then pass it to `job.addProfile()`

Python API: Adding Profiles

daxgen.py

```
1 # add a Profile Object
2 job.addProfile(Profile(Namespace.ENV, 'PATH', '/bin'))
3
4 # add a profile directly
5 job.profile(Namespace.CONDOR, "universe", "vanilla")
```

Call `job.profile()`, and set the profile directly.

Python API: Adding Profiles

workflow.py

```
1  from Pegasus.api import *
2
3  j = Job("office-hours")
4
5  # option 1: ensures valid key is entered, can set multiple profiles
6  j.add_pegasus_profile(clusters_num=1, stageout_clusters=1)
7  j.add_env(PATH="/bin", JAVA_HOME="/java/home")
8  j.add_condor_profile(universe="vanilla")
9  j.add_dagman_profile(retry=3)
10 j.add_globus_profile(count=1)
11 j.add_selector_profile(execution_site="condorpool")
12
13 # option 2: lower level, can set multiple profiles
14 j.add_profiles(Namespace.ENV, PATH="/bin", JAVA_HOME="/java/home")
15
16 # option 3: lower level, use when keys have non-alphanumeric characters
17 j.add_profiles(Namespace.CONDOR, key="+KeyName", value="val")
```

Python API: Adding Profiles

workflow.py

```
1  from Pegasus.api import *
2
3  j = Job("office-hours")
4
5  # option 1: ensures valid key is entered, can set multiple profiles
6  j.add_pegasus_profile(clusters_num=1, stageout_clusters=1)
7  j.add_env(PATH="/bin", JAVA_HOME="/java/home")
8  j.add_condor_profile(universe="vanilla")
9  j.add_dagman_profile(retry=3)
10 j.add_globus_profile(count=1)
11 j.add_selector_profile(execution_site="condorpool")
12
13 # option 2: lower level, can set multiple profiles
14 j.add_profiles(Namespace.ENV, PATH="/bin", JAVA_HOME="/java/home")
15
16 # option 3: lower level, use when keys have non-alphanumeric characters
17 j.add_profiles(Namespace.CONDOR, key="+KeyName", value="val")
```

Call `obj.add_<namespace>_profile(key1=val1, key2=val2, ...)`

Python API: Adding Profiles

workflow.py

```
1  from Pegasus.api import *
2
3  j = Job("office-hours")
4
5  # option 1: ensures valid key is entered, can set multiple profiles
6  j.add_pegasus_profile(clusters_num=1, stageout_clusters=1)
7  j.add_env(PATH="/bin", JAVA_HOME="/java/home")
8  j.add_condor_profile(universe="vanilla")
9  j.add_dagman_profile(retry=3)
10 j.add_globus_profile(count=1)
11 j.add_selector_profile(execution_site="condorpool")
12
13 # option 2: lower level, can set multiple profiles
14 j.add_profiles(Namespace.ENV, PATH="/bin", JAVA_HOME="/java/home")
15
16 # option 3: lower level, use when keys have non-alphanumeric characters
17 j.add_profiles(Namespace.CONDOR, key="+KeyName", value="val")
```

Call `obj.add_profiles(Namespace.<NS>, key1=val1, ...)`

Python API: Adding Profiles

workflow.py

```
1  from Pegasus.api import *
2
3  j = Job("office-hours")
4
5  # option 1: ensures valid key is entered, can set multiple profiles
6  j.add_pegasus_profile(clusters_num=1, stageout_clusters=1)
7  j.add_env(PATH="/bin", JAVA_HOME="/java/home")
8  j.add_condor_profile(universe="vanilla")
9  j.add_dagman_profile(retry=3)
10 j.add_globus_profile(count=1)
11 j.add_selector_profile(execution_site="condorpool")
12
13 # option 2: lower level, can set multiple profiles
14 j.add_profiles(Namespace.ENV, PATH="/bin", JAVA_HOME="/java/home")
15
16 # option 3: lower level, use when keys have non-alphanumeric characters
17 j.add_profiles(Namespace.CONDOR, key="+KeyName", value="val")
```

Call `obj.add_profiles(Namespace.<NS>, key="+-.something", value="val")`

Python API: Putting It All Together

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ## Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ## Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ## Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ## Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ## Plan and Run #####
41 wf.plan(submit=True).wait()
```

Python API: Putting It All Together

import Pegasus.api; setup

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

Python API: Putting It All Together

create properties file (config)

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

pegasus.properties

Python API: Putting It All Together

create replica catalog
(specify input files)

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

pegasus.properties

replicas.yml

Python API: Putting It All Together

create transformation catalog
(specify executables)

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

pegasus.properties

replicas.yml

transformations.yml

Python API: Putting It All Together

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

pegasus.properties

replicas.yml

transformations.yml

build abstract workflow

Python API: Putting It All Together

```
1 #!/usr/bin/env python3
2 import logging
3 from pathlib import Path
4
5 from Pegasus.api import *
6
7 logging.basicConfig(level=logging.INFO)
8
9 TOP_DIR = Path(__file__).parent.resolve()
10
11 ### Properties #####
12 props = Properties()
13 props["pegasus.mode"] = "development"
14 props.write()
15
16 ### Replica Catalog#####
17 rc = ReplicaCatalog()
18 rc.add_replica(site="local", lfn="input.txt", pfn=TOP_DIR / "input.txt")
19 rc.write()
20
21 ### Transformation Catalog#####
22 tc = TransformationCatalog()
23 compute = Transformation(
24     "compute.sh",
25     site="local",
26     pfn=TOP_DIR / "compute.sh",
27     is_stageable=True
28 )
29 tc.add_transformations(compute)
30 tc.write()
31
32 ### Workflow #####
33 wf = Workflow("office-hours-workflow")
34
35 in_file = File("input.txt")
36 out_file = File("output.txt")
37 j1 = Job(compute).add_args("OH").add_inputs(in_file).add_outputs(out_file)
38 wf.add_jobs(j1)
39
40 ### Plan and Run #####
41 wf.plan(submit=True).wait()
```

files will be consumed
by **pegasus-plan**

pegasus.properties

replicas.yml

transformations.yml

workflow.yml

plan and run workflow

Outline

- ~~Overview of Changes~~
- ~~Pegasus 4.X vs Pegasus 5.0 Workflow Scripts~~
- ~~Migrating to New Python API~~
- **Useful Tips**
- Resources

Useful Tips: Building Up the Replica Catalog

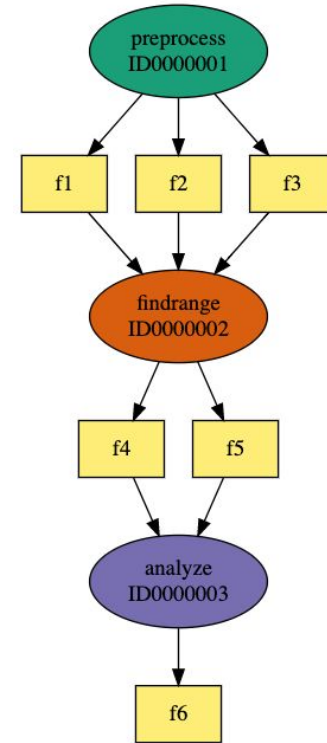
```
1 from pathlib import Path
2
3 from Pegasus.api import *
4
5 rc = ReplicaCatalog()
6
7 input_dir = Path(".").parent / "inputs"
8 for p in input_dir.iterdir():
9     rc.add_replica(site="local", lfn=p.name, pfn=p.resolve())
10
11 rc.write()
```

```
inputs
├── input_1.txt
├── input_2.txt
├── input_3.txt
├── input_4.txt
└── input_5.txt
```

Loop through input directory and add all files to the Replica Catalog object.
(also see [`pegasus-plan --input-dir dir1\[,dir2,...\]`](#))

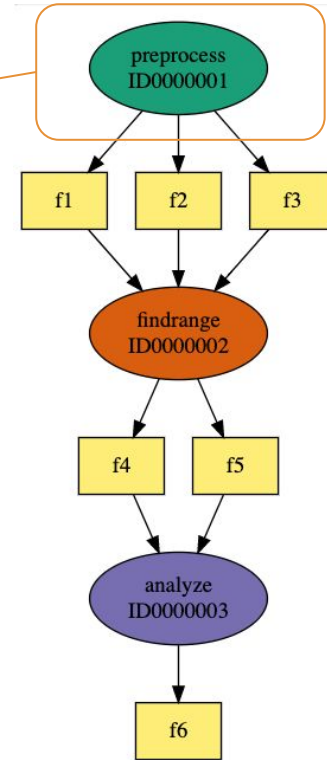
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



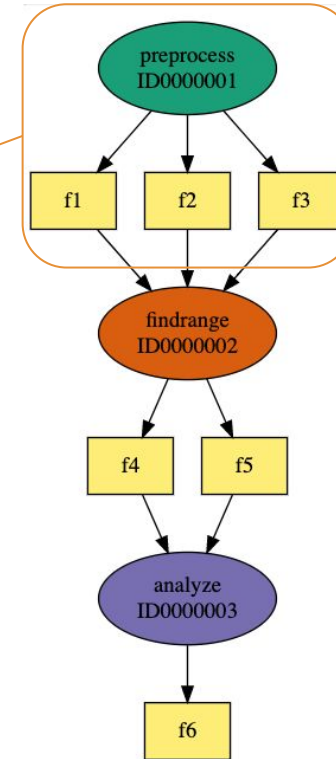
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



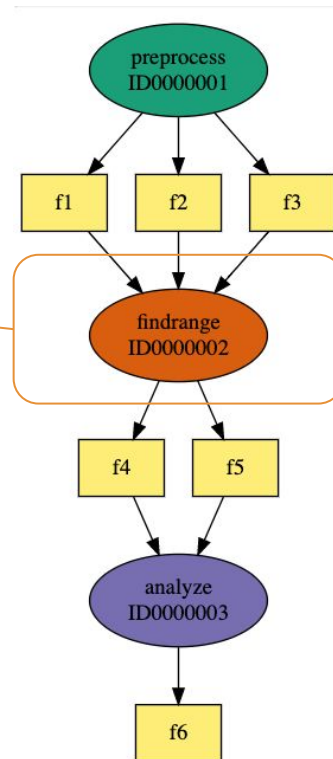
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



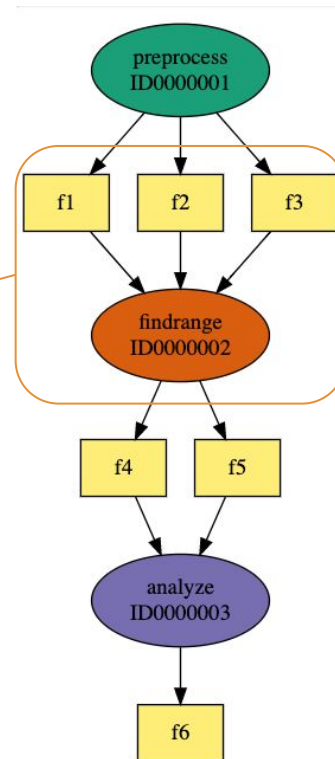
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



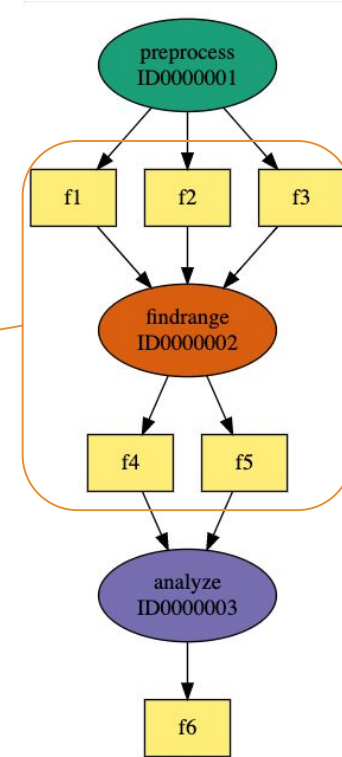
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



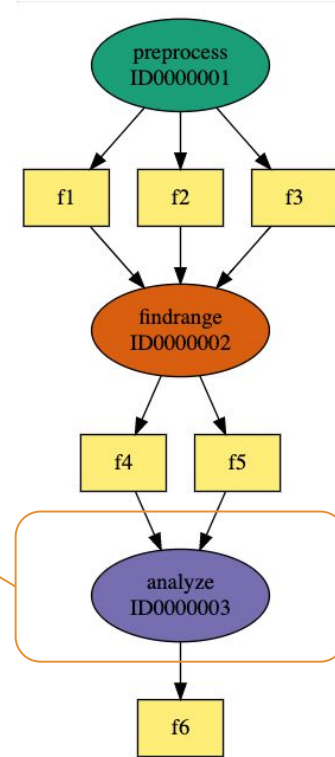
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



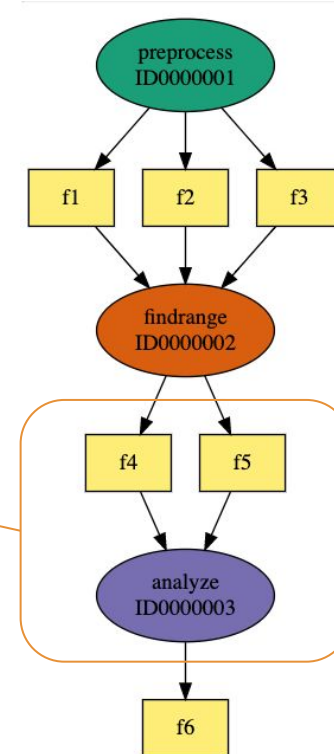
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



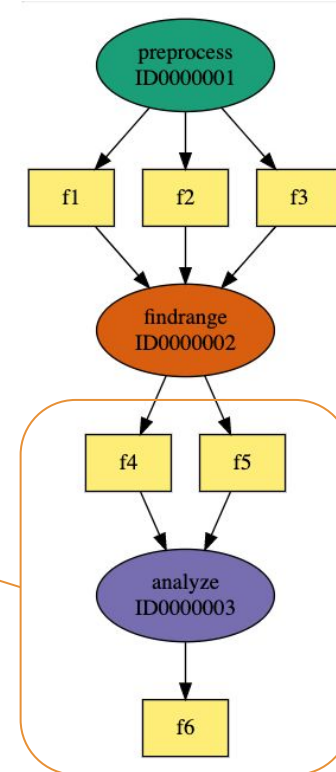
Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



Useful Tips: Chaining Jobs Together

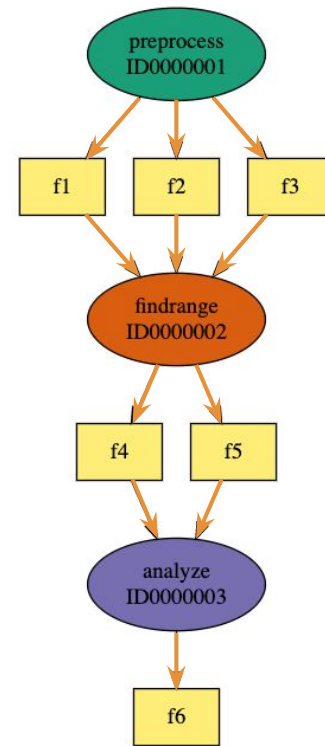
```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```



Useful Tips: Chaining Jobs Together

```
1 wf = Workflow("office-hours")
2 job_1 = Job("preprocess")
3 job_1.add_outputs(File("f1"), File("f2"), File("f3"))
4
5 job_2 = Job("findrange")
6 job_2.add_inputs(*job_1.get_outputs())
7 job_2.add_outputs(File("f4"), File("f5"))
8
9 job_3 = Job("analyze")
10 job_3.add_inputs(*job_2.get_outputs())
11 job_3.add_outputs(File("f6"))
12
13 wf.add_jobs(job_1, job_2, job_3)
```

By default, dependencies are inferred based on input/output files created/used by each job.

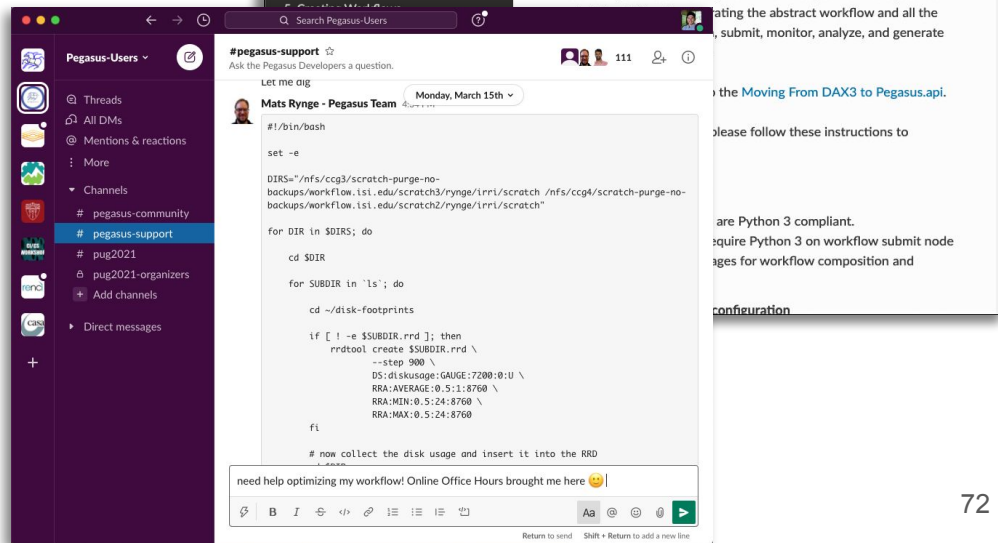
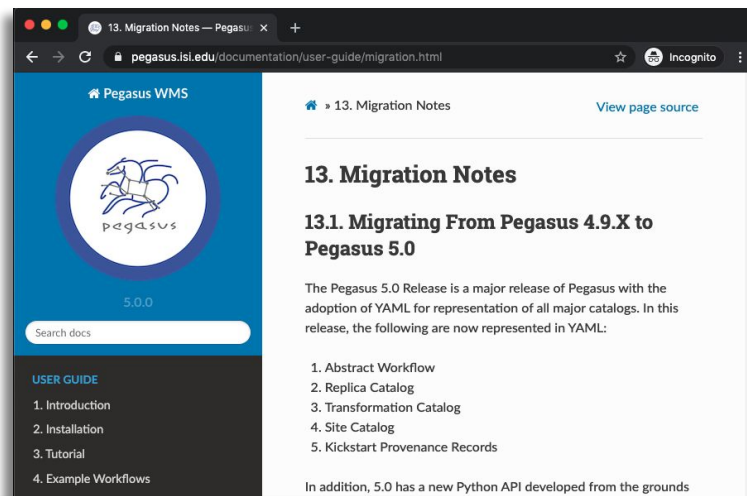


Outline

- ~~Overview of Changes~~
- ~~Pegasus 4.X vs Pegasus 5.0 Workflow Scripts~~
- ~~Migrating to New Python API~~
- ~~Useful Tips~~
- **Resources**

Resources

- Migration Notes
 - <https://pegasus.isi.edu/documentation/user-guide/migration.html>
- Python API Reference
 - <https://pegasus.isi.edu/documentation/reference-guide/api-reference.html>
- Pegasus Tutorial
 - <https://pegasus.isi.edu/documentation/user-guide/tutorial.html>
- Pegasus Team
 - <https://pegasus-users.slack.com/archives/C01KBPSM64S>





Pegasus est. 2001

Automate, recover, and debug scientific computations

Get Started

Pegasus Online Office Hours

<https://pegasus.isi.edu/blog/online-pegasus-office-hours>

Bi-monthly basis on the second Friday of the month, where we address user questions and also apprise the community of new developments.

Pegasus Website

<https://pegasus.isi.edu/>

Users Mailing List

pegasus-users@isi.edu

Pegasus Website

pegasus-support@isi.edu