
PipelineJobs System Documentation

Release 0.2.0

TACC Open Source

Feb 04, 2019

Components

1	Pipelines Manager	1
2	PipelineJobs Manager	9
3	PipelineJobs Indexer	13
4	PipelineJobs Agave Proxy	17
5	HTTP Authentication	25
6	Send a Bearer Token	27
7	Use a Nonce	29
8	Authorization Tokens	31
9	Reactor+App with PipelineJobs Integration	33
10	Indices and tables	35

CHAPTER 1

Pipelines Manager

This Reactor enables creation and management of Data Catalog **Pipelines**.

1.1 Pipelines

Pipelines are defined using document written in the **Pipeline JSON** schema. Here is a brief, slightly silly example:

```
{
  "name": "Tacobot 9001",
  "description": "Creates even better breakfast tacos out of silicon, ozone, and_
→shredded GPUs",
  "components": [
    {
      "id": "e1LqagMyqq4xB",
      "opts": {
        "cheese": true,
        "salsa": true,
        "eggs": false
      }
    },
    {
      "id": "bbq-brisket-0.5.0",
      "parameters": {
        "smoke": true,
        "quantity": 4
      }
    }
  ],
  "processing_levels": [
    "1"
  ],
  "accepts": [
    "CSV"
  ]
}
```

(continues on next page)

(continued from previous page)

```

    ],
    "produces": [
        "PNG"
    ]
}

```

A Pipeline is a human-readable name and description, a globally-unique string identifier, and a list of components that defines some number of Abaco Actors, Agave Apps, Deployed Containers, and Web Services. Additionally, one or more data “processing levels” are provided as well as the list of file types accepted and emitted. Of these fields, only `components` is used to create a Pipeline UUID that connects each Pipeline to its compute jobs.

To make this very tangible: In the example above, if *anything* in `components` changes, the result will have to be a new UUID. If the name, file types, description, or description change, the UUID remains the same.

1.2 Messages

All Pipeline management actions are accomplished by sending JSON-formatted messages to the **Pipelines Manager** Reactor.

1.3 Create a New Pipeline

A new pipeline can be registered by sending its JSON document to the **Pipelines Manager** as a message like so:

```

$ abaco run -m "$(jq -c . my_pipeline.json)" G1p783Pxp1BB

gOvQRGRVPP0zZ

# Wait a few seconds

$ abaco logs G1p783Pxp1BB gOvQRGRVPP0zZ

EGe6NKeo8Oy5 DEBUG Action selected: create
EGe6NKeo8Oy5 INFO Created pipeline 1064aaf1-459c-5e42-820d-b822aa4b3990 with update_
↪token 0df45d5e9e0f31e2

```

Note the pipeline’s UUID (1064aaf1-459c-5e42-820d-b822aa4b3990). This is needed to configure **PipelineJobs** that reference this **Pipeline**. Also note the update token (0df45d5e9e0f31e2). This is needed to update the Pipeline at a later date.

1.4 Update a Pipeline

A pipeline’s *components* cannot be updated, but its human-readable name and description can be. To accomplish this, edit the pipeline JSON document with new or amended values and send it to the **Pipelines Manager** along with a valid update token.

```

$ abaco run -m "$(jq -c . my_pipeline.json)" -q token=0df45d5e9e0f31e2 G1p783Pxp1BB

e5QKEW8L0BeZ4

```

(continues on next page)

(continued from previous page)

```
# Wait a few seconds

$ abaco logs G1p783Pxp1aBB e5QKEW8L0Bez4

EkBWRvV1gKG1a DEBUG Action selected: update
EkBWRvV1gKG1a INFO Updated pipeline 1064aaf1-459c-5e42-820d-b822aa4b3990 with update_
↪token 0df45d5e9e0f31e2
```

1.5 Retire a Pipeline

A pipeline cannot be deleted, but it can be retired from active service.

Coming soon...

1.6 JSON Schemas

Listing 1: pipeline_manager_create

```
1 {
2     "$schema": "http://json-schema.org/draft-07/schema#",
3     "$id": "https://schema.catalog.sd2e.org/schemas/pipeline_manager_create.json",
4     "title": "PipelinesDefinition",
5     "description": "A Pipeline record. POST to PipelinesManager to create.",
6     "definitions": {
7         "container_repo": {
8             "type": "string",
9             "description": "a Linux container image repository and tag"
10        },
11        "agave_app": {
12            "description": "Agave application",
13            "type": "object",
14            "properties": {
15                "id": {
16                    "type": "string",
17                    "description": "the distinct 'app.id' for the_
↪Agave app"
18                },
19                "inputs": {
20                    "type": "object",
21                    "description": "predefined inputs for jobs_
↪spawned by the app"
22                },
23                "parameters": {
24                    "type": "object",
25                    "description": "predefined parameters for the_
↪jobs spawned by the app",
26                    "properties": {
27                        "CONTAINER_IMAGE": {
28                            "$ref": "#/definitions/
↪container_repo",
29                            "description": "Linux_
↪container image repository and tag for the deployed app"
```

(continues on next page)

(continued from previous page)

```

30         }
31     }
32 }
33 },
34     "required": ["id", "inputs", "parameters"]
35 },
36     "container": {
37         "description": "Deployed Linux container",
38         "type": "object",
39         "properties": {
40             "repo": {
41                 "$ref": "#/definitions/container_repo",
42                 "description": "Linux container image_
↪ repository and tag for the service"
43             },
44             "hash": {
45                 "type": "string",
46                 "description": "Linux container image hash"
47             },
48             "options": {
49                 "type": "object",
50                 "description": "Deployment options"
51             }
52         },
53         "required": ["repo"]
54     },
55     "service": {
56         "description": "External networked resource",
57         "type": "object",
58         "properties": {
59             "identifier": {
60                 "type": "string",
61                 "description": "An identifier for a specific_
↪ instance of a service"
62             },
63             "uri": {
64                 "type": "string",
65                 "description": "Canonical URI for the resource
↪ ",
66                 "format": "uri"
67             },
68             "options": {
69                 "type": "object",
70                 "description": "Additional descriptive_
↪ attributes"
71             }
72         },
73         "required": ["uri"]
74     },
75     "abaco_actor": {
76         "description": "Abaco Reactor",
77         "type": "object",
78         "properties": {
79             "id": {
80                 "type": "string",
81                 "description": "Distinct 'actor.id' for the_
↪ Reactor"

```

(continues on next page)

(continued from previous page)

```

82         },
83         "repo": {
84             "$ref": "#/definitions/container_repo",
85             "description": "Linux container image_
↪ repository and tag for the deployed actor"
86         },
87         "options": {
88             "type": "object",
89             "description": "Predefined runtime options_
↪ for the Reactor"
90         },
91     },
92     "required": ["id", "repo"]
93 },
94 "pipeline_type": {
95     "description": "the general class of action done by the_
↪ pipeline",
96     "type": "string",
97     "enum": ["generic-process", "data-transfer", "metadata-
↪ management", "primary-etl", "secondary-etl"],
98     "default": "primary-etl"
99 },
100 "processing_level": {
101     "description": "a data processing level",
102     "type": "string",
103     "enum": ["0", "1", "2", "3"]
104 },
105 "collections_level": {
106     "description": "a data processing level",
107     "type": "string",
108     "enum": ["reference", "user_file", "challenge_problem",
↪ "experiment", "sample", "measurement", "file", "pipeline", "job", "product"]
109 },
110 },
111 "type": "object",
112 "properties": {
113     "name": {
114         "type": "string"
115     },
116     "description": {
117         "type": "string"
118     },
119     "components": {
120         "description": "an unordered array of apps and actors in the_
↪ pipeline (required)",
121         "type": "array",
122         "items": {
123             "anyOf": [{
124                 "$ref": "#/definitions/agave_app"
125             },
126             {
127                 "$ref": "#/definitions/abaco_actor"
128             },
129             {
130                 "$ref": "#/definitions/service"
131             },
132             {

```

(continues on next page)

(continued from previous page)

```

133         "$ref": "#/definitions/container"
134     }
135 ]
136 }
137 },
138 "collections_levels": {
139     "type": "array",
140     "items": {
141         "$ref": "#/definitions/collections_level"
142     },
143     "description": "level(s) of data input that the pipeline acts_
↪ upon",
144     "default": []
145 },
146 "processing_levels": {
147     "type": "array",
148     "description": "level(s) of data product produced by the_
↪ pipeline",
149     "items": {
150         "$ref": "#/definitions/processing_level"
151     },
152     "default": []
153 },
154 "pipeline_type": {
155     "items": {
156         "$ref": "#/definitions/pipeline_type"
157     }
158 },
159 "accepts": {
160     "type": "array",
161     "description": "file types accepted by the pipeline",
162     "items": {
163         "type": "string"
164     },
165     "default": ["*"]
166 },
167 "produces": {
168     "type": "array",
169     "description": "file types produced by the pipeline",
170     "items": {
171         "type": "string"
172     },
173     "default": ["*"]
174 },
175 "__options": {
176     "description": "private object for passing runtime options to_
↪ a pipeline (optional)",
177     "type": "object"
178 }
179 },
180 "required": ["name", "components"],
181 "additionalProperties": false
182 }

```

Listing 2: pipeline_manager_update

```
1 {
2     "$schema": "http://json-schema.org/draft-07/schema#",
3     "$id": "https://schema.catalog.sd2e.org/schemas/pipeline_manager_update.json",
4     "title": "PipelinesDefinitionUpdate",
5     "description": "Update a Pipeline record",
6     "type": "object",
7     "properties": {
8         "uuid": {
9             "description": "The job UUID",
10            "type": "string"
11        },
12        "body": {"type": "object"},
13        "action": {
14            "type": "string",
15            "enum": ["update"]
16        },
17        "token": {
18            "description": "an authorization token issued when the
19↪pipeline was created",
20            "type": "string",
21            "minLength": 16,
22            "maxLength": 17
23        },
24        "__options": {
25            "type": "object",
26            "description": "an object used to pass runtime options to a
27↪reactor (private, optional)"
28        }
29    },
30    "required": ["uuid", "action", "body", "token"],
31    "additionalProperties": false
32 }
```

PipelineJobs Manager

This Reactor manages updates to **PipelineJobs** once they are created by other processes using the **ManagedPipelineJob** and **ReactorManagedPipelineJob** classes.

2.1 Update Job State

PipelineJobs Manager can update a job's state via three mechanisms:

1. Receipt of a JSON-formatted **pipelinejob_manager_event**
2. Receipt of URL parameters sufficient to form a **pipelinejob_manager_event**
3. Receipt of a **agave_job_callback** POST combined with **uuid**, **status** and **token** URL parameters

These named document formats are documented below in *JSONSchemas*.

2.1.1 JSON Event Messages

The default method for updating a **PipelineJob**'s state is to send a JSON message to the Manager actor. Here is an example of a **finish** event for job `1073f4ff-c2b9-5190-bd9a-e6a406d9796a`, which will indicate the job has completed primary computation and archiving steps.

```
{
  "uuid": "1073f4ff-c2b9-5190-bd9a-e6a406d9796a",
  "name": "finish",
  "token": "0dc73dc3ff39b49a"
}
```

This can be sent directly over HTTP like so:

```
curl -XPOST -H "Authorization: Bearer 969dl1396c43b0b810387e4da840cb37" \
  --data '{"uuid": "1073f4ff-c2b9-5190-bd9a-e6a406d9796a", \
  "token": "0dc73dc3ff39b49a", \
```

(continues on next page)

(continued from previous page)

```
"name": "finish"}' \
https://api.tacc.cloud/actors/v2/<actorId>/messages
```

It can also be sent from within a Recator like so:

```
rx = Reactor()
manager_id = '<actorId>'
finish_mes = { 'uuid': '1073f4ff-c2b9-5190-bd9a-e6a406d9796a',
               'name': 'finish',
               'token': '0dc73dc3ff39b49a' }
rx.send_message(manager_id, finish_mes)
```

2.1.2 URL Parameters Event

The `uuid`, and `event`, `token` fields can be sent as URL parameters in an HTTP POST. The contents of the POST body will be attached to the event if one is present. Our **finish** event expressed as URL parameters looks like:

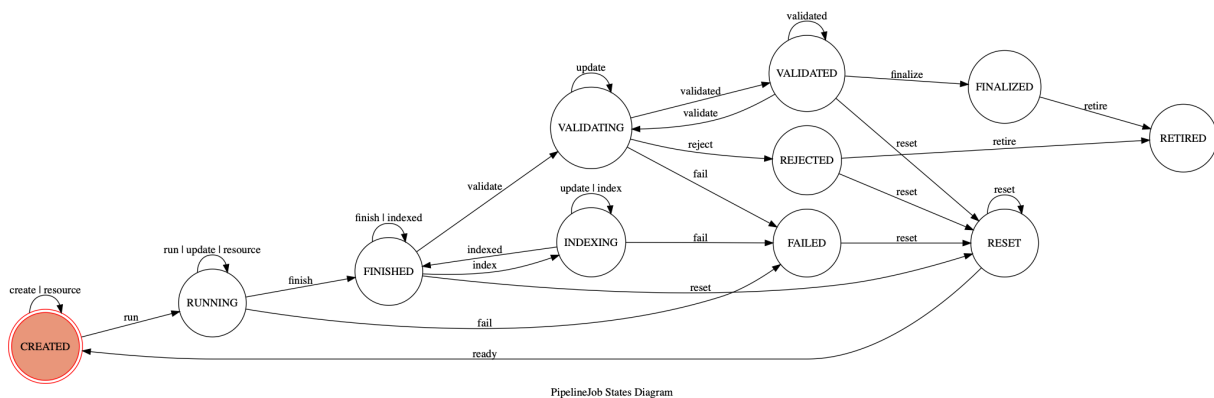
```
curl -XPOST --data '{"arbitrary": "key value data"}' \
https://api.tacc.cloud/actors/v2/<actorId>/messages?uuid=1073f4ff-c2b9-5190-bd9a-
e6a406d9796a&\
event=finish&token=0dc73dc3ff39b49a
```

2.1.3 Agave Jobs Notification

HTTP POST body and URL parameters are combined to link the Agave Jobs system with **PipelineJobs**, which is quite handy as Agave jobs often are enlisted to do the computational heavy lifting in analysis workflows. This approach is demonstrated in the **demo-jobs-reactor-app** repository.

2.2 PipelineJobs Events

The state of every PipelineJob proceeds through a defined lifecycle, where transitions occur in response to receipt of named events. This is illustrated in the following image:



PipelineJobs Manager accepts any of the events (lower case words attached to the graph edges) save for **create**, which is reserved for other agents.

2.3 Authentication

POSTs to a **PipelineJobs Manager** must be authenticated by one of two means:

1. Send a valid TACC.cloud OAuth2 Bearer token with the request
2. Include a special URL parameter called a **nonce** with the HTTP request

2.4 JSON Schemas

Listing 1: agave_job_callback

```

1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "$id": "https://schema.catalog.sd2e.org/schemas/agave_job_callback.json",
4   "title": "AgaveJobsPost",
5   "description": "A job status POST from the Agave API jobs service",
6   "type": "object",
7   "properties": {
8     "owner": {
9       "type": "string"
10    },
11    "appId": {
12      "type": "string"
13    },
14    "executionSystem": {
15      "type": "string"
16    },
17    "archivePath": {
18      "type": "string"
19    },
20    "archiveSystem": {
21      "type": "string"
22    },
23    "status": {
24      "type": "string"
25    },
26    "inputs": {},
27    "parameters": {},
28    "_links": {}
29  },
30  "required": [
31    "appId",
32    "executionSystem",
33    "owner",
34    "archivePath",
35    "archiveSystem",
36    "status",
37    "inputs",
38    "parameters",
39    "_links"
40  ]
41 }
```

Listing 2: pipelinejob_manager_event

```
1 {
2     "$schema": "http://json-schema.org/draft-07/schema#",
3     "$id": "https://schema.catalog.sd2e.org/schemas/pipelinejob_manager_event.json",
4     "title": "PipelineJobManagerEvent",
5     "description": "A state-change event for a PipelineJob",
6     "type": "object",
7     "properties": {
8         "uuid": {
9             "$ref": "pipelinejob_uuid.json"
10        },
11        "name": {
12            "$ref": "pipelinejob_eventname.json"
13        },
14        "data": {
15            "description": "Additional information to attach to the event",
16            "type": "object",
17        },
18        "token": {
19            "$ref": "update_token.json"
20        }
21    },
22    "required": ["uuid", "name"],
23    "additionalProperties": false
24 }
```

PipelineJobs Indexer

This Reactor indexes the contents of **ManagedPipelineJob** archive paths. It implements two actions: **index** and **indexed**. It is normally run automatically via **PipelineJobs Manager** when a job enters the **FINISHED** state, but can also be activated on its own.

3.1 Index a Job

PipelineJobs Indexer can receive an **index** request via:

1. A JSON-formatted **pipelinejob_index** document
2. URL parameters that replicate a **pipelinejob_index** document

Here are the critical fields to request indexing:

1. **uuid** ID for job to be indexed (must validate as a known job)
2. **name** This is always **index**
3. **token** The job's update token (optional for now)
4. **level** The processing level for output files (default: **1**)
5. **filters** List of **url-encoded** Python regex that select a subset of **archive_path**

3.1.1 Index Request as JSON

This message will index outputs of job 1079f67e-0ef6-52fe-b4e9-d77875573860 as level “2” products, sub-selecting only files matching **sample\..uw_biofab\..141715** and **sample-uw_biofab-141715**.

```
{
  "uuid": "1079f67e-0ef6-52fe-b4e9-d77875573860",
  "name": "index",
  "filters": [
```

(continues on next page)

(continued from previous page)

```
    "sample%5C.uw_biofab%5C.141715",
    "sample-uw_biofab-141715"
  ],
  "level": "2",
  "token": "0dc73dc3ff39b49a"
}
```

3.1.2 Index Request as URL Params

```
curl -XPOST \
  https://<tenantUrl>/actors/v2/<actorId>/messages?uuid=1073f4ff-c2b9-5190-bd9a-
  ↪e6a406d9796a&\
    level=2&token=0dc73dc3ff39b49a&name=index --data '{"filters": ["sample.uw_biofab.
  ↪141715", "sample-uw_biofab-141715"]}'
```

Note: Remember that `filters` cannot currently be passed as URL parameters.

3.2 Mark as Indexed

PipelineJobs Indexer can receive an **index** request via JSON message or URL parameters. Here is an example.

```
{
  "uuid": "1079f67e-0ef6-52fe-b4e9-d77875573860",
  "name": "indexed",
  "token": "0dc73dc3ff39b49a"
}
```

PipelineJobs Indexer sends itself an **indexed** message after completing an indexing action. Thus, it is not usually necessary to send one manually and in fact, should be avoided. Documentation on the **indexed** event is included here mostly for the sake of completeness.

3.3 Authentication

POSTs to a **PipelineJobs Indexer** must be authenticated by one of two means:

1. Send a valid TACC.cloud OAuth2 Bearer token with the request
2. Include a special URL parameter called a **nonce** with the HTTP request

3.4 JSON Schemas

Listing 1: pipelinejob_index

```
1 {
2   "$schema": "http://json-schema.org/draft-07/schema#",
3   "$id": "https://schema.catalog.sd2e.org/schemas/pipelinejob_index.json",
```

(continues on next page)

(continued from previous page)

```

4     "title": "PipelineJobIndexEvent",
5     "description": "Request indexing of a completed PipelineJob's archive path",
6     "type": "object",
7     "properties": {
8         "uuid": {
9             "$ref": "pipelinejob_uuid.json"
10        },
11        "name": {
12            "type": "string",
13            "enum": [
14                "index"
15            ]
16        },
17        "filters": {
18            "type": "array",
19            "description": "List of Python regular expressions defining which output
↳ files to associate with the job. Omit entirely if you do not want to apply
↳ filtering.",
20            "items": {
21                "type": "string"
22            }
23        },
24        "level": {
25            "$ref": "processing_level.json"
26        },
27        "token": {
28            "$ref": "update_token.json"
29        }
30    },
31    "required": ["uuid", "name"],
32    "additionalProperties": false
33 }

```

Listing 2: pipelinejob_indexed

```

1 {
2     "$schema": "http://json-schema.org/draft-07/schema#",
3     "$id": "https://schema.catalog.sd2e.org/schemas/pipelinejob_indexed.json",
4     "title": "PipelineJobIndexedEvent",
5     "description": "Mark PipelineJob indexing as completed.",
6     "type": "object",
7     "properties": {
8         "uuid": {
9             "$ref": "pipelinejob_uuid.json"
10        },
11        "name": {
12            "type": "string",
13            "enum": [
14                "indexed"
15            ]
16        },
17        "token": {
18            "$ref": "update_token.json"
19        }
20    },
21    "required": ["uuid", "name"],

```

(continues on next page)

(continued from previous page)

```
22     "additionalProperties": false
23 }
```

This Reactor provides a generalized proxy for running Agave API jobs such that their inputs, parameterization, and outputs are connected to (and thus discoverable from within) the Data Catalog.

4.1 Register Agave App as a Pipeline

Before an Agave App can be run by this proxy, three things must happen:

1. It must be architected to fit the PipelineJobs workflow
2. It must be public or shared with user **sd2eadm**
3. It must be registered as a Data Catalog Pipeline

4.1.1 App Architecture

The app must generate filenames that are distinguishable between runs. This is enforced to prevent accidentally overwriting of files when multiple jobs share an archiving destination. Furthermore, the app definition and any interior runtime logic must use fully-qualified Agave files URLs to define inputs. Finally, the app's `id` must be unique not only in the **Agave Apps Catalog** (this is automatically enforced) but also in the **Data Catalog Pipelines** collection.

4.1.2 Share or Publish the App

Coming soon...

4.1.3 Registering a Pipeline

Coming soon...

4.2 Launching a Managed Agave Job

Construct and send a message including the following components to the **PipelineJobs Agave Proxy** Reactor.

1. An Agave job definition
2. A metadata linkage parameter
3. Optional control parameters

Note: The `agave_pipelinejob` format is documented in *JSONSchemas*.

4.2.1 Agave Job Definition

The Agave job definition must be included as a subdocument in the message. To illustrate this, start with a basic Agave job definition: Here is an example for an imaginary Agave app `tacobot9000-0.1.0u1`.

```
{
  "appId": "tacobot9000-0.1.0u1",
  "name": "TACObot job",
  "inputs": { "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt" },
  "parameters": { "salsa": true, "avocado": false, "cheese": true },
  "maxRunTime": "01:00:00"
}
```

To launch this via the Agave, this document would be sent directly to the `/apps` endpoint. To send it instead to the proxy, move it to key `job_definition` in a JSON document.

```
{
  "job_definition": {
    "appId": "tacobot9000-0.1.0u1",
    "name": "TACObot job",
    "inputs": {
      "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
    },
    "parameters": {
      "salsa": true,
      "avocado": false,
      "cheese": true
    },
    "maxRunTime": "01:00:00"
  }
}
```

4.2.2 Metadata Linkage Parameter

An explicit linkage to objects in the Data Catalog must be established. This is done via the `parameters` key, which must contain a valid value for one of the following:

- `experiment_id`
- `sample_id`
- `measurement_id`

Either single values or an array of values may be passed, and either the readable text value may be provided or the corresponding UUID.

Which Parameter to Pass

A PipelineJob is always linked to a set of measurements by way of the linkage parameter. The job's archive path is also determined by the linkage parameter. To illustrate:

If a job's `measurement_id=['measurement.tacc.1234', 'measurement.tacc.2345']`, it will be linked to these two measurements and its archive path will end with a hash of the two `measurement_id` values.

Assuming those measurements are children of `sample.tacc.abcde` and the only linkage parameter sent was `sample_id='sample.tacc.abcdef'`, the job will still be linked to all the child measurements of that sample. Its archive path will end with a hash of `sample.tacc.abcde`. However, if both `measurement_id` and `sample_id` are passed, the linkages are made to the specified measurement(s) while the archive path is a function of the `sample_id` value(s).

For `experiment_id`, the specific samples are linked to the job and the archive path is a function of `experiment_id` value(s).

This design allows files generated by the job to be linked to only one level of the metadata hierarchy, while allowing collection of outputs at higher levels of organization in the file system.

Here is a worked example of the current example job request, as it stands:

```
{
  "parameters": {
    "sample_id": "sample.tacc.abcde"
  },
  "job_definition": {
    "appId": "tacobot9000-0.1.0u1",
    "name": "TACObot job",
    "inputs": {
      "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
    },
    "parameters": {
      "salsa": true,
      "avocado": false,
      "cheese": true
    },
    "maxRunTime": "01:00:00"
  }
}
```

4.2.3 Additional Control Parameters

Job behavior can be refined with additional control parameters.

instanced

Each PipelineJob has a distinct archive path derived from its Pipeline UUID, the `data` dictionary passed at `job init()` and/or `setup()`, and a function of its linkage parameters to experiments, samples, or measurements. To avoid inadvertent over-writes, the archive path is extended with an *instancing directory* named in the form `adjective-animal-YYYYMMDDTHHmssZ`. To avoid use of the instancing directory, include `instanced: false` in the job request message.

Example: "instanced": false

index_patterns

The default behavior of the PipelineJobs System is to index every file found under a job's archive path to be linked to that specific job. To subselect only specific files, it is possible to include one or more Python regular expressions in `index_patterns`. Only files matching these patterns will be linked to the job.

Example: "index_patterns": []

processing_level

The default behavior of the PipelineJobs System is to index files under a job's archive path as processing level "1". To change this, an alternative `processing_level` may be passed in the job request message.

Example: "processing_level": "2"

Note: Only one automatic indexing configuration can be active for a given job. Additional indexing actions with other configurations may be initiated by sending a message directly to **PipelineJobs Indexer**

4.3 Job Life Cycle

Here is complete record from the Pipelines system showing how the information from job creation and subsequent events is stored and discoverable. A few key highlights:

- The top-level data field holds the original parameterization of the job
- Three events are noted in the `history`: create, run, finish
- The actor and execution for the managing instance of **PipelineJobs Agave Proxy** are available under `agent` and `task`, respectively

```
1 {
2   "agent": "https://api.tacc.cloud/actors/v2/G46vjoAVzGkkz",
3   "archive_path": "/products/v2/103f877a7ab857d182807b75af4eab6e/
4   ↳106bd127e2d257acb9belled06042e68/eligible-awk-20181127T173243Z",
5   "archive_system": "data-sd2e-community",
6   "data": {
7     "appId": "urrutia-novel_chassis_app-0.1.0",
8     "archivePath": "",
9     "inputs": {
10      "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
11    },
12    "maxRunTime": "01:00:00",
13    "name": "TACObot job",
14    "parameters": {
15      "avocado": false,
16      "cheese": true,
17      "salsa": true
18    }
19  },
20  "derived_from": [
    "1022efa3-4480-538f-a581-f1810fb4e0c3"
```

(continues on next page)

(continued from previous page)

```

21 ],
22 "generated_by": [
23     "106bd127-e2d2-57ac-b9be-11ed06042e68"
24 ],
25 "history": [
26     {
27         "data": {
28             "appId": "tacobot9000-0.1.0u1",
29             "inputs": {
30                 "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
31             },
32             "maxRunTime": "01:00:00",
33             "name": "TACObot job",
34             "parameters": {
35                 "avocado": false,
36                 "cheese": true,
37                 "salsa": true
38             }
39         },
40         "date": "2018-12-08T00:08:32.000+0000",
41         "name": "create"
42     },
43     {
44         "data": {
45             "appId": "tacobot9000-0.1.0u1",
46             "archive": true,
47             "archivePath": "/products/v2/103f877a7ab857d182807b75af4eab6e/
↪106bd127e2d257acb9belled06042e68/eligible-awk-20181127T173243Z",
48             "archiveSystem": "data-tacc-cloud",
49             "batchQueue": "normal",
50             "created": "2018-12-07T18:08:37.000-06:00",
51             "endTime": null,
52             "executionSystem": "hpc-tacc-stampede2",
53             "id": "7381691026605150696-242ac11b-0001-007",
54             "inputs": {
55                 "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
56             },
57             "lastUpdated": "2018-12-07T18:09:40.000-06:00",
58             "maxRunTime": "01:00:00",
59             "memoryPerNode": 1,
60             "name": "TACObot job",
61             "nodeCount": 1,
62             "outputPath": "tacobot/job-7381691026605150696-242ac11b-0001-007-
↪TACObot-job",
63             "owner": "tacobot",
64             "parameters": {
65                 "avocado": false,
66                 "cheese": true,
67                 "salsa": true
68             },
69             "processorsPerNode": 1,
70             "startTime": null,
71             "status": "RUNNING",
72             "submitTime": "2018-12-07T18:09:40.000-06:00"
73         },
74         "date": "2018-12-08T00:10:12.000+0000",
75         "name": "run"

```

(continues on next page)

(continued from previous page)

```

76     },
77     {
78         "data": {
79             "appId": "tacobot9000-0.1.0u1",
80             "archive": true,
81             "archivePath": "/products/v2/103f877a7ab857d182807b75af4eab6e/
↪106bd127e2d257acb9belled06042e68/eligible-awk-20181127T173243Z",
82             "archiveSystem": "data-tacc-cloud",
83             "batchQueue": "normal",
84             "created": "2018-12-07T18:08:37.000-06:00",
85             "endTime": null,
86             "executionSystem": "hpc-tacc-stampede2",
87             "id": "7381691026605150696-242ac11b-0001-007",
88             "inputs": {
89                 "file1": "agave://data.tacc.cloud/examples/tacobot/test1.txt"
90             },
91             "lastUpdated": "2018-12-07T18:53:20.000-06:00",
92             "maxRunTime": "01:00:00",
93             "memoryPerNode": 1,
94             "name": "TACObot job",
95             "nodeCount": 1,
96             "outputPath": "tacobot/job-7381691026605150696-242ac11b-0001-007-
↪TACObot-job",
97             "owner": "tacobot",
98             "parameters": {
99                 "avocado": false,
100                 "cheese": true,
101                 "salsa": true
102             },
103             "processorsPerNode": 1,
104             "startTime": "2018-12-07T18:09:49.000-06:00",
105             "status": "FINISHED",
106             "submitTime": "2018-12-07T18:09:40.000-06:00"
107         },
108         "date": "2018-12-08T00:53:45.000+0000",
109         "name": "finish"
110     }
111 ],
112 "last_event": "finish",
113 "pipeline_uuid": "106bd127-e2d2-57ac-b9be-11ed06042e68",
114 "session": "casual-bass",
115 "state": "FINISHED",
116 "task": "https://api.tacc.cloud/actors/v2/G46vjoAVzGkkz/executions/Myp6wvklV0zgQ",
117 "updated": "2018-12-08T00:53:45.000+0000",
118 "uuid": "10743f9e-f5ae-5b4c-859e-6774ef4ab08b"
119 }

```

4.4 JSON Schemas

Listing 1: agave_pipelinejob

```

1 {
2     "$schema": "http://json-schema.org/draft-07/schema#",
3     "$id": "https://schema.catalog.sd2e.org/schemas/agave_pipelinejob.json",

```

(continues on next page)

(continued from previous page)

```

4  "title": "AgavePipelineJob",
5  "description": "Launch an Agave job as a PipelineJob",
6  "type": "object",
7  "definitions": {
8      "datacatalog_field": {
9          "type": "object",
10         "anyOf": [
11             {
12                 "required": [
13                     "sample_id"
14                 ],
15             },
16             {
17                 "required": [
18                     "experiment_design_id"
19                 ],
20             },
21             {
22                 "required": [
23                     "experiment_id"
24                 ],
25             },
26             {
27                 "required": [
28                     "measurement_id"
29                 ],
30             }
31         ],
32         "properties": {
33             "sample_id": {
34                 "type": "string",
35                 "pattern": "^sample.(uw_biofab|transcriptic|ginkgo|emerald).$"
36             },
37             "experiment_design_id": {
38                 "type": "string",
39                 "$ref": "experiment_reference.json"
40             },
41             "experiment_id": {
42                 "type": "string",
43                 "pattern": "^experiment.(uw_biofab|transcriptic|ginkgo|emerald).$"
44             },
45             "measurement_id": {
46                 "type": "string",
47                 "pattern": "^measurement.(uw_biofab|transcriptic|ginkgo|emerald).$"
48             }
49         }
50     },
51 },
52 "properties": {
53     "parameters": {
54         "$ref": "#/definitions/datacatalog_field"
55     },
56     "job_definition": {
57         "description": "An Agave API job definition",
58         "type": "object"
59     },
60     "archive_path": {

```

(continues on next page)

(continued from previous page)

```

61         "description": "Optional Agave URN defining the job's archive path",
62         "$ref": "agave_files_uri.json"
63     },
64     "instanced": {
65         "description": "Whether the generated archive path should be instanced_
↪with a randomized session",
66         "type": "boolean",
67         "value": true
68     },
69     "data": {
70         "description": "Optional dict-like object describing the job's run-time_
↪parameterization",
71         "type": "object"
72     },
73     "index_patterns": {
74         "type": "array",
75         "description": "List of Python regular expressions defining which output_
↪files to associate with the job. Omit entirely if you do not want to apply_
↪filtering.",
76         "items": {
77             "type": "string"
78         }
79     },
80     "processing_level": {
81         "description": "Defaults to '1' if not provided",
82         "$ref": "processing_level.json"
83     }
84 },
85 "required": [
86     "job_definition",
87     "parameters"
88 ],
89 "additionalProperties": false
90 }

```

HTTP Authentication

All POSTs to PipelineJobs components must be authenticated. There are two mechanisms by which this can be accomplished.

1. Send a valid TACC.cloud OAuth2 Bearer token with the request
2. Include a special URL parameter called a **nonce** with the HTTP request

Send a Bearer Token

Requests authenticated using a Bearer token run as the TACC.cloud identity of the user that issued the token. This is usually your account (or a role account if you have appropriate credentials).

6.1 Usage Example

```
curl -XPOST -H "Authorization: Bearer 969d11396c43b0b810387e4da840cb37" \  
  --data '{"uuid": "1073f4ff-c2b9-5190-bd9a-e6a406d9796a", \  
  "token": "0dc73dc3ff39b49a", \  
  "name": "finish"}' \  
  https://<tenantUrl>/actors/v2/<actorId>/messages
```


CHAPTER 7

Use a Nonce

```
curl -XPOST --data '{"arbitrary": "key value data"}' \
  https://<tenantUrl>/actors/v2/<actorId>/messages?uuid=1073f4ff-c2b9-5190-bd9a-
↪e6a406d9796a&\
  name=finish&token=0dc73dc3ff39b49a&\
  x-nonce=TACC_XXXXXXXxYz
```

Authorization Tokens

Update and delete actions for Pipelines and PipelineJobs are authorized via an additional token scoped to the Pipelines system. There are two types:

- Document update token
- Administrative action token

A document update token authorizes management of **one specific** Pipeline or Job. The token for a given document is returned after each update or management action and allows future actions to be taken at any time in the future. An administrative action token authorizes **any action** for **any Pipeline or Job**. They are obtained via an external process, and only by persons who possessing the **Administrator Token API Key**. Because they are powerful, administrative action tokens expire after 30 seconds.

8.1 Sending a Token

A token must be sent with the requested action or event. It can be included in the request message as field ("token": "<token>") or included as a URL parameter token=<token>.

Reactor+App with PipelineJobs Integration

This is an example implementation of a Reactor that launches an Agave app, where both Reactor and the App integrate with the PipelineJobs system to update the state and history of a PipelineJob.

9.1 Key Aspects

Note the following as you go through the code and configuration of this project:

1. Use of the `ManagedPipelineJob` class to instantiate the job
2. Addition of Agave job notifications allowing the job to send updates
3. Configuration of this Reactor via the `mongodb` and `pipelines` stanzas in `config.yml`
4. Use of Abaco API keys, or `_nonces_`, in `secrets.json`
5. Leveraging shared Abaco Reactors

9.2 Deploying your own Reactor

Note: Familiarity with configuration and deployment of Abaco Reactor is assumed.

Obtain the following values from the Infrastructure team for *secrets.json*:

- `_REACTOR_MONGODB_AUTHN`
- `_REACTOR_PIPELINES_JOB_MANAGER_NONCE`

Confirm the following values with the Infrastructure team for *config.yml*:

- `pipelines.job_manager_id`
- `pipelines.job_indexer_id`

Danger: Never include values from `secrets.json` in `config.yml` file. If you do this, they will be committed forever to a publicly-accessible Docker image when the Reactor is deployed.

Once you've set up the configuration and secrets, build then deploy your Reactor.

9.3 Send a Message

```
abaco run -m '{"data":{"key1":"value 1"}, "sample_id": "sample.tacc.20001}" <ACTOR_ID>
```

CHAPTER 10

Indices and tables

- `genindex`
- `modindex`
- `search`