iris hep
Institute for Research & Innovation
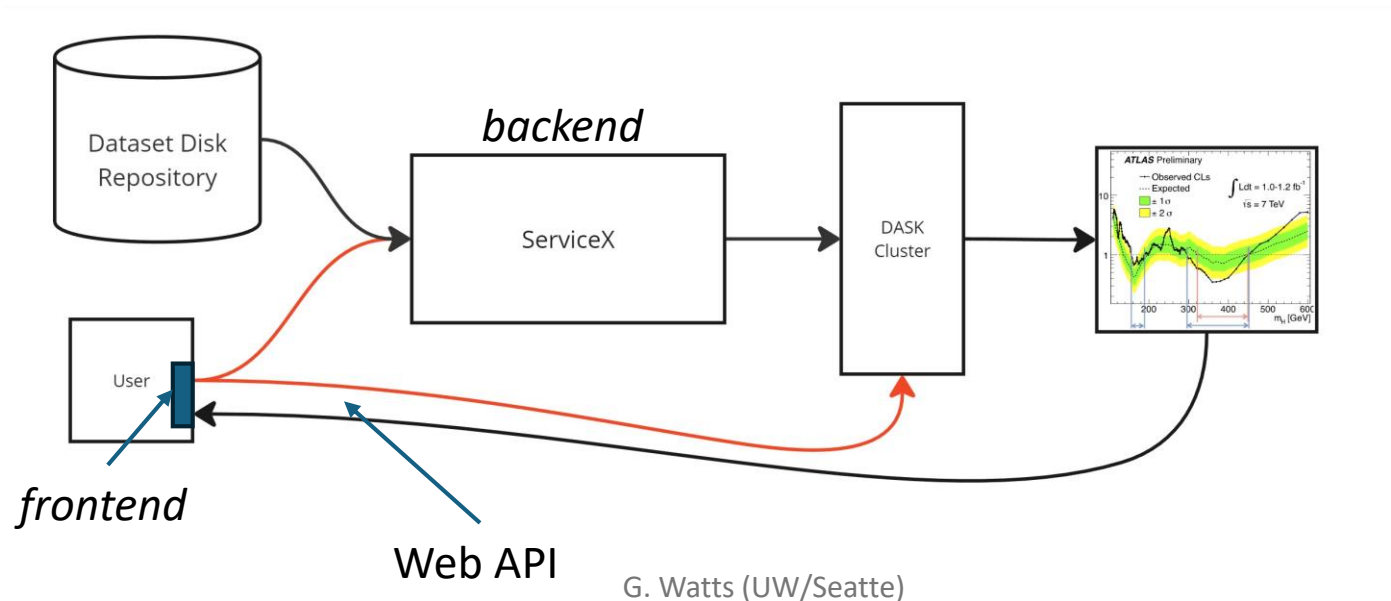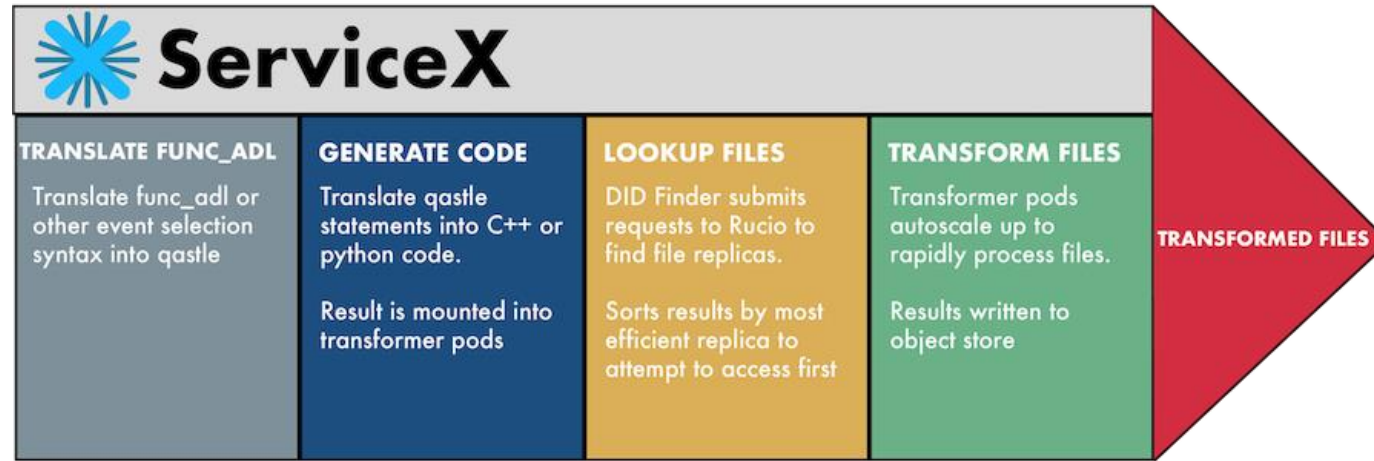in Software for High Energy Physics

# 3.0

G. Watts (UW/Seattle)

For the ServiceX Team

Nov 1, 2024

# ServiceX Architecture Reminder

# Recent Developments in ServiceX

- Backend Developments (ServiceX itself)
  - 1.4.1 (August 7th) -> 1.5.X (Sept 20th)
  - Lessons from the IDAP 200 Gpbs test: reduce lost internal scaling messages, small updates to how we transform the data, lots of small stability improvements.
  - New transformers/codegenerators to run plain-old-python
  - The WebAPI did not change, however!
  - Won't really discuss further
- Front End (Library to enable user interaction)
  - 3.0 released Sept 20th (after 5 months of development)
  - Big (breaking) change in how the user interacts with ServiceX
  - Pulled several ideas from the community into the central library (e.g. datasets)
  - Queries can be coded with typed classes, dictionaries, or yaml text files now!
  - DOCUMENTATION!!

# Getting Started: Help!

Find our new documentation on *readthedocs*:

ServiceX 3.0.0 documentation

Missing something? Spot an error?

Create an issue

Or… the source for the documentation is in the repo

Submit a MR

# How does it work?

Data Source (FileList, Rucio, etc.)
- Where to get the data from

Query
- How to transform the data
- Raw, func_adl, etc.

```python
from servicex import Sample, ServiceXSpec, query, dataset, deliver


spec = ServiceXSpec(
    Sample=[
        Sample(
            Name="UprootRaw_Typed",
            Dataset=dataset.FileList(
                [
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                ]
            ),
            Query=query.UprootRaw(
                [
                    {
                        "treename": "CollectionTree",
                        "filter_name": "AnalysisElectronsAuxDyn.pt",
                    }
                ]
            )
        )
    ]
)

print(f"Files: {deliver(spec)}")
```

G. Watts (UW/Seatte)

# How does it work?

Sample
- Consists of the data set and query together
- Produces a set of files
- Labeled with the Name

```python
from servicex import Sample, ServiceXSpec, query, dataset, deliver


spec = ServiceXSpec(
    Sample=[
        Sample(
            Name="UprootRaw_Typed",
            Dataset=dataset.FileList(
                [
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                ]
            ),
            Query=query.UprootRaw(
                [
                    {

                        "treename": "CollectionTree",
                        "filter_name": "AnalysisElectronsAuxDyn.pt",

                    }
                ]
            )
        )
    ]
)

print(f"Files: {deliver(spec)}")
```

# How does it work?

Samples
- You can submit 1 or 100 samples at a time
- Output will always be a dict of everything

All of these objects take lots of extra parameter that allow you to specify explicitly:
- Code generator
- Backend ServiceX location
- Transformer image, etc.

The servicex.yaml file is still required!!

```python
from servicex import Sample, ServiceXSpec, query, dataset, deliver


spec = ServiceXSpec(
    Sample=[
        Sample(
            Name="UprootRaw_Typed",
            Dataset=dataset.FileList(
                [
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                ]
            ),
            Query=query.UprootRaw(
                [
                    {
                        "treename": "CollectionTree",
                        "filter_name": "AnalysisElectronsAuxDyn.pt",
                    }
                ]
            )
        )
    ]
)

print(f"Files: {deliver(spec)}")
```

# How does it work?

deliver changes the spec with the Samples into files.

- There is an async version of this coming soon!
- It will handle all interaction with the ServiceX backend via the WebAPI

```python
from servicex import Sample, ServiceXSpec, query, dataset, deliver


spec = ServiceXSpec(
    Sample=[
        Sample(
            Name="UprootRaw_Typed",
            Dataset=dataset.FileList(
                [
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSL
                ]
            ),
            Query=query.UprootRaw(
                [
                    {
                        "treename": "CollectionTree",
                        "filter_name": "AnalysisElectronsAuxDyn.pt",
                    }
                ]
            )
        )
    ]
)

print(f"Files: {deliver(spec)}")
```

# Using a Dictionary instead…

We expect:

- Dictionary is the easiest to use for quick oneoffs
- Typed classes will be used by libraries and frameworks
- The YAML is attractive because it can be checked into git directly!

At its core, the frontend only uses the typed classes – everything is translated into that!

```python
from servicex import query, dataset, deliver


spec = {
    'Sample': [{
        'Name': "UprootRaw_Dict",
        'Dataset': dataset.FileList(
            [
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.
            ]
        ),
        'Query': query.UprootRaw(
            [
                {
                    "treename": "CollectionTree",
                    "filter_name": "AnalysisElectronsAuxDyn.pt",
                }
            ]
        )
    }]
}

print(f"Files: {deliver(spec)}")
```

# Using YAML

The YAML file itself:

```yaml
Sample:
  - Name: UprootRaw_YAML
    Dataset: !FileList
      [
        "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.37019878
        "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.37019878
        "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE.37019878
      ]
    Query: !UprootRaw |
      [{"treename":"CollectionTree", "filter_name": "AnalysisElectronsAuxDyn.pt"}]
```

Use in SX:

```python
from servicex import deliver

print(
    deliver("config_Uproot_FuncADL.yaml")
)
```

There is not yet a full text representation for all query languages!

# Other Queries: Func_ADL xAOD (Fully Typed)

Func_adl

ServiceX Spec and Delivery
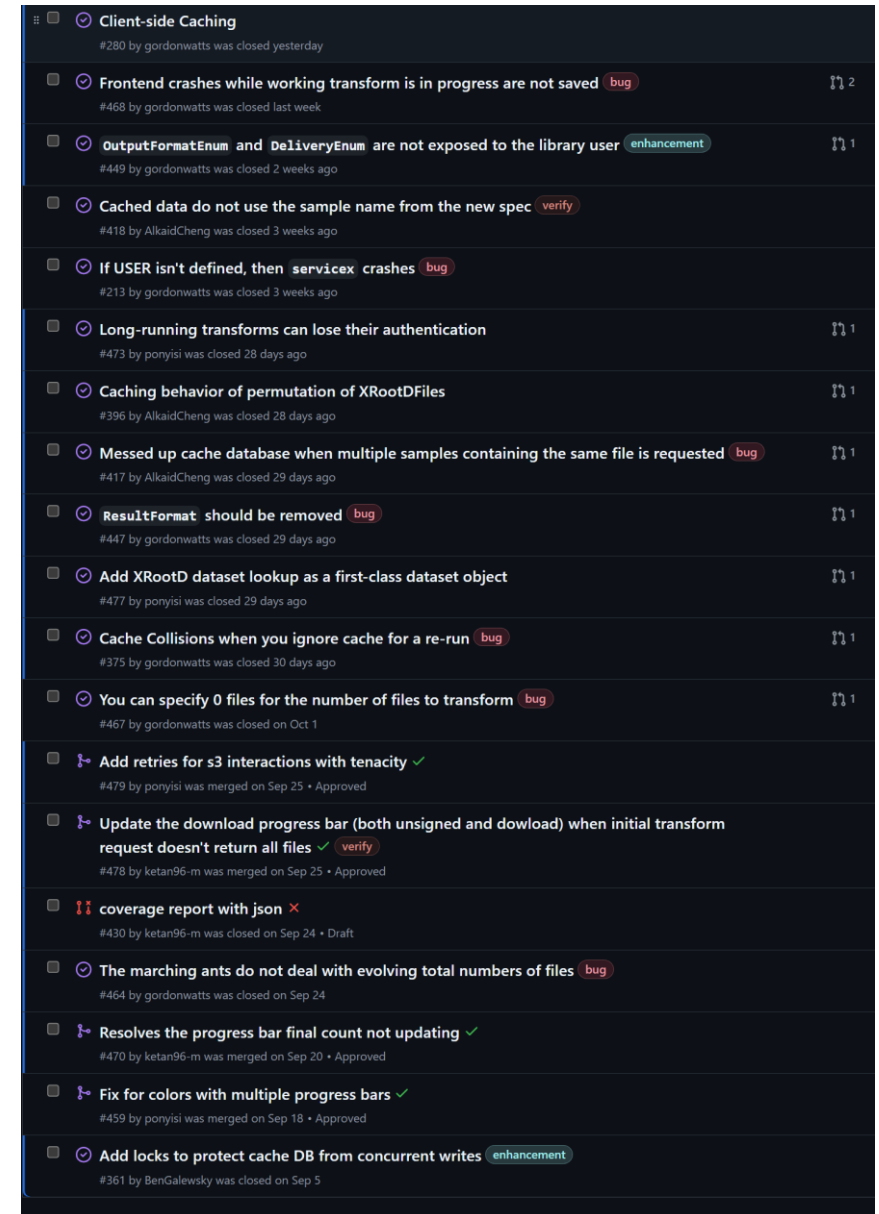
```python
from servicex import query as q, deliver, dataset


def func_adl_xaod_simple():
    query = q.FuncADL_ATLASr22()  # type: ignore
    jets_per_event = query.Select(lambda e: e.Jets('AnalysisJets'))
    jet_info_per_event = jets_per_event.Select(
        lambda jets: {
            'pt': jets.Select(lambda j: j.pt()),
            'eta': jets.Select(lambda j: j.eta())
        }
    )

    spec = {
        'Sample': [{
            'Name': "func_adl_xAOD_simple",
            'Dataset': dataset.FileList(
                [
                    "root://eospublic.cern.ch//eos/opendata/atlas/rucio/mc20_13TeV/DAOD_PHYSLIT
                ]
            ),
            'Query': jet_info_per_event
        }]
    }
    files = deliver(spec, servicex_name="servicex-uc-af")
    assert files is not None, "No files returned from deliver! Internal error"
    return files
```

G. Watts (UW/Seatte)

# Status

- The Analysis Grand Challenge and the IDAP 200 Gbps code has been converted to use 3.0

- Being used in the wild as well (been out about 1.5 months).

- However, there are some rough edges
  - 3.0.1's branch already has a significant number of updates
  - Will release "soon"

- We also have a backlog of new features where we are planning for 3.1

# Backup

# Full Python Function Transformer!!

Python Function

This allows one to extract non-ntuple-like objects from the file (e.g. cutflow, etc.)

ServiceX Spec and Delivery

```python
from servicex import query, dataset, deliver


def run_query(input_filenames=None):
    import uproot  # type: ignore
    with uproot.open({input_filenames: "CollectionTree"}) as o:
        br = o.arrays("AnalysisElectronsAuxDyn.pt")
    return br



spec = {
    'Sample': [{
        'Name': "PythonFunction_Dict",
        'Dataset': dataset.FileList(
            [
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE
                "root://eospublic.cern.ch//eos/opendata/atlas/rucio/data16_13TeV/DAOD_PHYSLITE
            ]
        ),
        'Query': query.PythonFunction().with_uroot_function(run_query)
    }]
}


print(f"Files: {deliver(spec)}")
```