

Capturing and recording provenance information from distributed workflows

Dr. Nelis Drost

nelis.drost@auckland.ac.nz

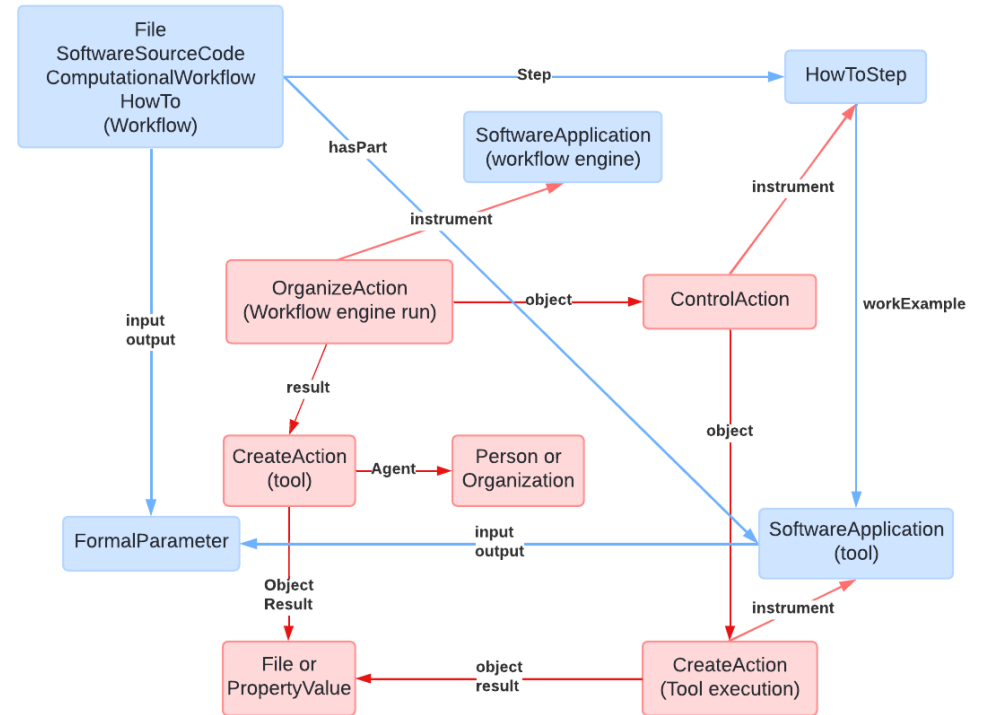
Augustus Ellerm

Prof. Mark Gahagen

Prof. Benjamin Adams

Dr. Ryan Chard

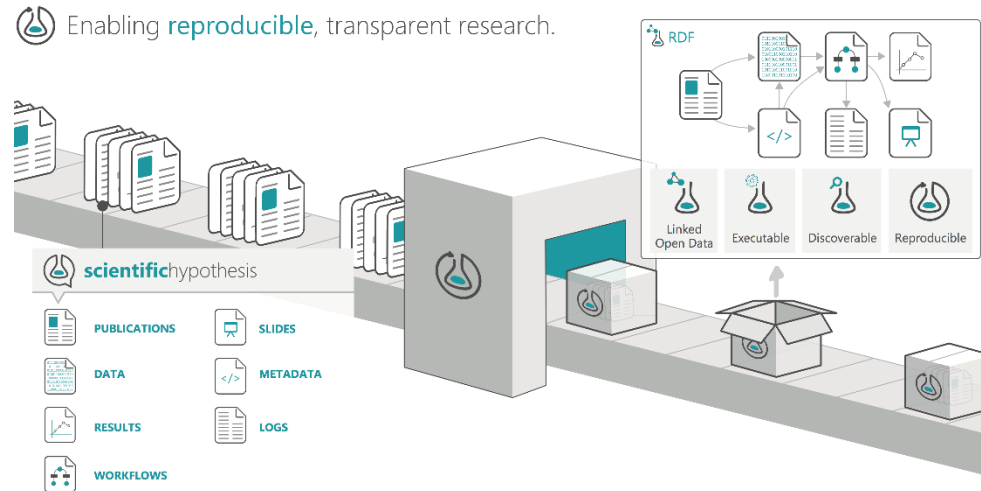
Prof. Ian Foster



What is data provenance?

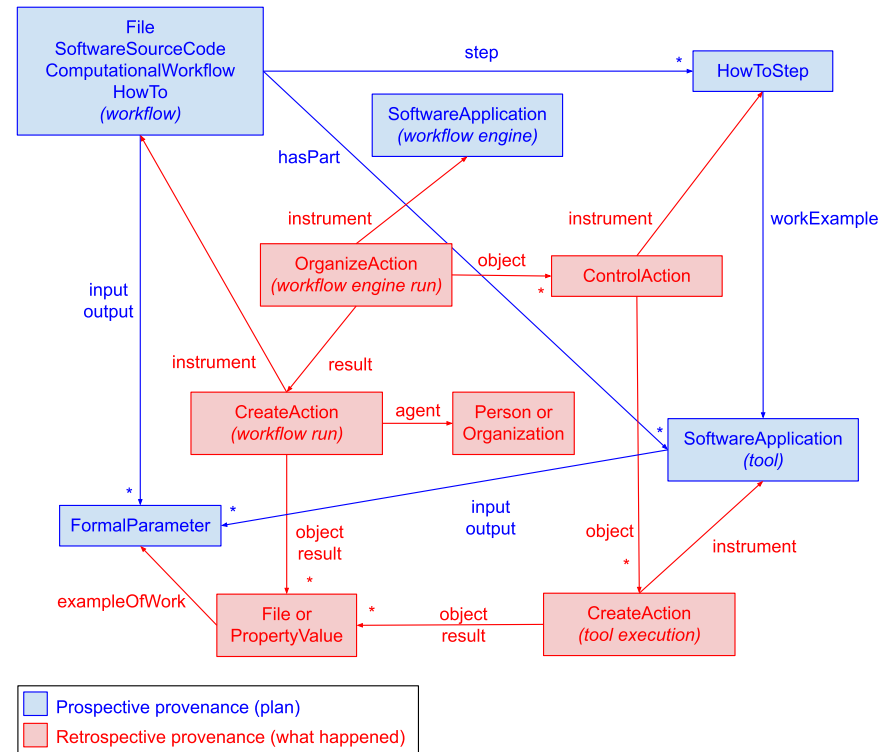
Metadata about how data was generated:

- Data sources: DOIs, UUIDs, etc.
- Authors, organisations, funders
- Instrument IDs
- Code repositories
- Configuration files
- Process



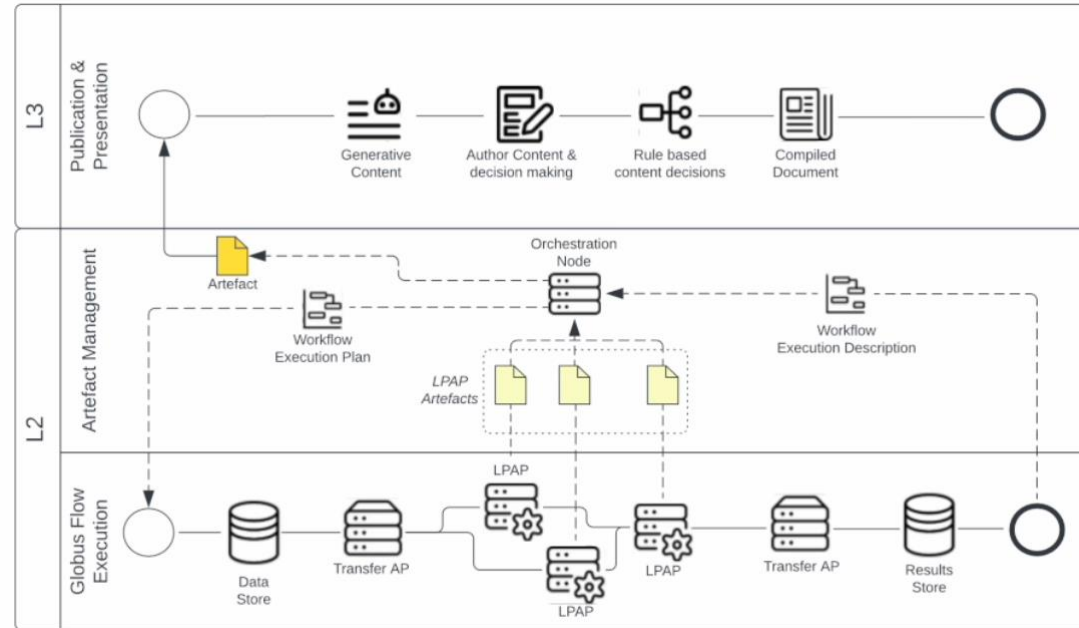
The RO-Crate Provenance Run Crate

- json-ld format
- Extends schema.org + bioschemas
- Details a workflow for generating data/results
- Describes both the planned (*prospective*) and actual (*retrospective*) steps taken



Live Publication

- Machine reproducible paper
- Automatically updateable
- Discovery of computational workflows
- Reuse and remixing of analyses



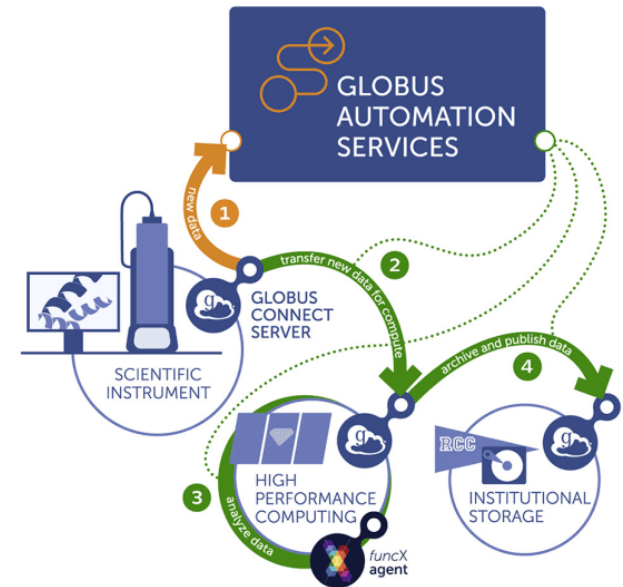
Integration with workflow management tools

- Workflow tools automate a series of computational steps
- Offer improved reproducibility
- May or may not produce provenance metadata
- Some tooling exists to produce RO-Crates from workflow scripts

| Implementation | Profile | Version URL/DOI | Example |
|----------------------------|-----------------|--|---|
| runcrate | Provenance | runcrate 0.5.0 or later | 10.5281/zenodo.7774351 |
| Galaxy | Workflow | Galaxy 23.1.1 or later | 10.5281/zenodo.7785861 |
| COMPSS | Workflow | compss 3.2 or later | 10.5281/zenodo.7788030 |
| StreamFlow | Provenance | Streamflow 0.2.0.dev10 | 10.5281/zenodo.7911906 |
| WfExS | Workflow | WfExS 0.10.1 or later | 10.5281/zenodo.10091550 |
| Sapporo | Workflow | sapporo-service 1.5.1 or later | 10.5281/zenodo.10134581 |
| Autosubmit | Workflow | Autosubmit v4.0.100 or later | 10.5281/zenodo.8144612 |
| Nextflow | Provenance | (nf-prov in development) | example |

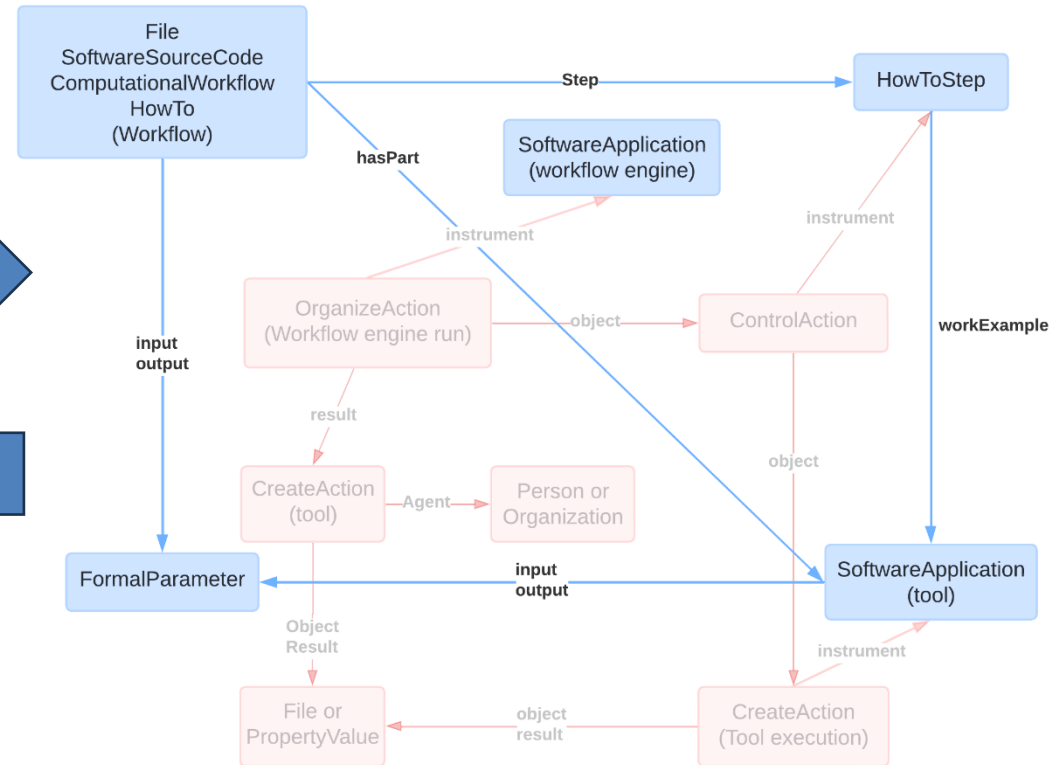
Provenance from federated systems

- Globus provides automated data transfer
- Globus Compute allows for automated computational flows
- The Globus ecosystem potentially enables federated workflows – combining data and computation from disparate sources
- We aim to produce provenance information describing such flows and making them reproducible



The challenge – lossless recording of workflow execution plan

```
"RevTxt": {
  "Comment": null,
  "Type": "Action",
  "ActionUrl": "https://compute.actions.globus.org",
  "ExceptionOnActionFailure": false,
  "Parameters": {
    "tasks": [
      {
        "endpoint.$": "$.input.compute_endpoint",
        "function.$": "$.input.rev_txt_function_id",
        "payload.$": "$.input.RevTxt"
      }
    ]
  },
  "ResultPath": "$.RevTxt",
  "WaitTime": 300,
  "Next": "Transfer_provenance_rev_txt"
},
```



Building on Gladier – capturing formal parameters

- Gladier simplifies the creation of computational workflows
- We can extend it to provide the additional information needed
- i.e.: Formal Parameters

```
58
59
60 1 usage
61 def sort_txt(input_file: str, output_file: str, reverse: bool):
62     with open(input_file, 'r') as f:
63         lines = f.readlines()
64         lines = [line.strip() for line in lines]
65         lines = [line for line in lines if line]
66         lines.sort()
67         if reverse:
68             lines = lines[::-1]
69
70     from pathlib import Path
71     Path(output_file).parent.mkdir(parents=True, exist_ok=True)
72     with open(output_file, 'w') as f:
73         f.write('\n'.join(lines))
```

Drost, 13/05/2024 12:52 pm • E

Building on Gladier – capturing formal parameters

- Gladier simplifies the creation of computational workflows
- We can extend it to provide the additional information needed
- i.e.: Formal Parameters

- Lesson: be more opinionated

```
58
59
60 1 usage
61 def sort_txt(input_file: str, output_file: str, reverse: bool):
62     with open(input_file, 'r') as f:
63         lines = f.readlines()
64         lines = [line.strip() for line in lines]
65         lines = [line for line in lines if line]
66         lines.sort()
67         if reverse:
68             lines = lines[::-1]
69
70     from pathlib import Path
71     Path(output_file).parent.mkdir(parents=True, exist_ok=True)
72     with open(output_file, 'w') as f:
73         f.write('\n'.join(lines))
```

Drost, 13/05/2024 12:52 pm • E

Distributed Step Crates

- Capture run time, *retrospective*, information about:
 - Tool execution
 - Hardware/software environment
- Record file creation
- Returned to the orchestration server between steps

```
"Transfer_provenance_rev_txt": {
  "Comment": "Transfer a file or directory in Globus",
  "Type": "Action",
  "ActionUrl": "https://actions.automate.globus.org/transfer/transfer",
  "Parameters": {
    "source_endpoint_id.$": "$.input.prov_compute_GCS_id",
    "destination_endpoint_id.$": "$.input.orchestration_server_endpoint_id",
    "transfer_items": [
      Drost, 13/05/2024 12:52 pm • Example flow that replicates th
      {
        "recursive": true,
        "source_path.=": "`$.RevTxt.details.results[0].task_id` + '.crate'",
        "destination_path.=": "`$.input._provenance_crate_destination_directory` + '/' + `$.RevT
      }
    ]
  },
  "ResultPath": "$.Transfer_provenance_rev_txt",
  "WaitTime": 600,
  "Next": "TransferRT_ST"
```

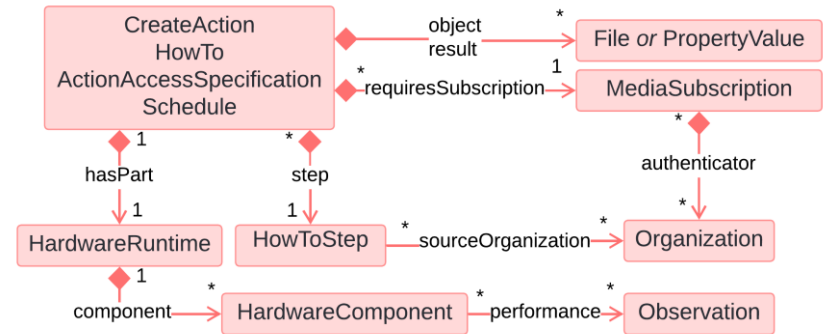
Distributed Step Crates

- Capture run time, *retrospective*, information about:
 - Tool execution
 - Hardware/software environment
- Record file creation
- Returned to the orchestration server between steps
- Lesson: none?

```
"Transfer_provenance_rev_txt": {
  "Comment": "Transfer a file or directory in Globus",
  "Type": "Action",
  "ActionUrl": "https://actions.automate.globus.org/transfer/transfer",
  "Parameters": {
    "source_endpoint_id.$": "$.input.prov_compute_GCS_id",
    "destination_endpoint_id.$": "$.input.orchestration_server_endpoint_id",
    "transfer_items": [
      Drost, 13/05/2024 12:52 pm • Example flow that replicates th
      {
        "recursive": true,
        "source_path.=": "`$.RevTxt.details.results[0].task_id` + '.crate'",
        "destination_path.=": "`$.input._provenance_crate_destination_directory` + '/' + `$.RevT
      }
    ]
  },
  "ResultPath": "$.Transfer_provenance_rev_txt",
  "WaitTime": 600,
  "Next": "TransferRT_ST"
```

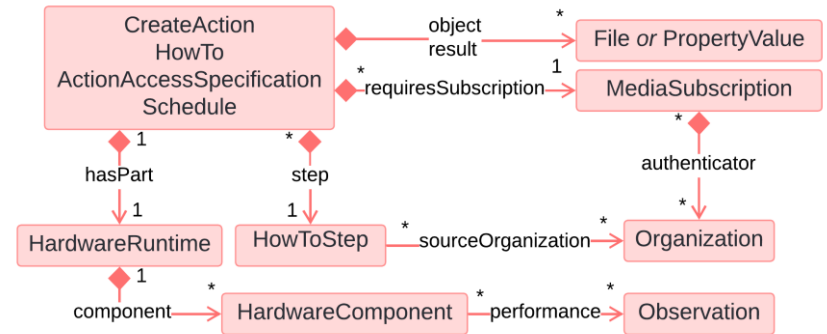
Crate extensions

- Hardware specs
 - Captured via a custom compute node
 - Subset captured by wrapped compute function
 - Returned as part of the Distributed Step Crate
- Access/authentication information



Crate extensions

- Hardware specs
 - Captured via a custom compute node
 - Subset captured by wrapped compute function
 - Returned as part of the Distributed Step Crate
- Access/authentication information
- Lesson: don't shoehorn everything into schema.org



Checking everything works

- Checking RO-Crates by hand is tedious and error prone
- Diff tools don't handle linked json
- *Retrospective* types may not be present – including links
- TDD works well for iterative development

```
"@graph": {
  {
    "@id": "WEP_json#RevTxt"
  },
  {
    "@id": "#f6174665-ae7d-4b03-b7c3-36716a311614",
    "@type": "ParameterConnection",
    "sourceParameter": {
      "@id": "main#$.input.RevTxt.input_file"
    },
    "targetParameter": {
      "@id": "$.input.RevTxt.input_file"
    }
  },
  {
    "@id": "WEP_json#RevTxt",
    "@type": "SoftwareApplication",
    "description": null,
    "input": [
      {
        "@id": "WEP_json#$.input.RevTxt.input_file"
      },
      {
        "@id": "WEP_json#$.input.RevTxt.output_file"
      }
    ],
    "name": "RevTxt",
    "output": [
      {
        "@id": "WEP_json#$.input.RevTxt.output_file"
      }
    ]
  }
},
```

```
{ "ecotext": [
  "https://w3id.org/ro/crate/1.1/context",
  "https://w3id.org/terms/workflow-run/context"
] },
{
  "@graph": [
    {
      "@id": "ro-crate-metadata.json",
      "@type": "CreativeWork",
      "about": [ { "id": "/" } ],
      "conformsTo": [
        { "id": "https://w3id.org/ro/crate/1.1" },
        { "id": "https://w3id.org/workflowhub/workflow-ro-crate/1.0" }
      ]
    },
    {
      "@id": "/",
      "@type": "Dataset",
      "conformsTo": [
        { "id": "https://w3id.org/rdfun/process/0.1" },
        { "id": "https://w3id.org/ro/rdfun/workflow/0.1" },
        { "id": "https://w3id.org/rdfun/provenance/0.1" },
        { "id": "https://w3id.org/workflowhub/workflow-ro-crate/1.0" }
      ],
      "hasPart": [
        { "id": "packed.cwl" },
        { "id": "327c7aedff6b9942a7c0b080dc5a7ff61376" },
        { "id": "b9214668cc45331b2c2282b772a5063dbd284" },
        { "id": "97fe1b50b4532e6c708537964b0e2e3e363aa3f" }
      ],
      "mainEntity": [ { "id": "packed.cwl" } ],
      "mentions": [
        { "id": "8a154d4d3-0bcc-4e35-bb8f-a2d6cd7dc49" }
      ]
    },
    {
      "@id": "https://w3id.org/rdfun/process/0.1",
      "@type": "CreativeWork",
      "name": "Process Run Crate",
      "version": "0.1"
    },
    {
      "@id": "https://w3id.org/rdfun/workflow/0.1",
      "@type": "CreativeWork",
      "name": "Workflow Run Crate",
      "version": "0.1"
    },
    {
      "@id": "https://w3id.org/rdfun/provenance/0.1",
      "@type": "CreativeWork",
      "name": "Provenance Run Crate",
      "version": "0.1"
    },
    {
      "@id": "https://w3id.org/workflowhub/workflow-ro-crate/1.0",
      "@type": "CreativeWork",
      "name": "Workflow RO-Crate",
      "version": "1.0"
    },
    {
      "@id": "packed.cwl",
      "@type": [ "File", "SoftwareSourceCode", "ComputationalWorkflow", "HowTo" ],
      "hasPart": [
        { "id": "packed.cwl#revtool.cwl" },

```


Lessons learned

- Start with the scaffolding:
 - Validation
 - Glazier
- Complexity is unpredictable
- Be opinionated
- Follow standards where possible, but not at all costs
- Start smaller

Any Questions?

- nelis.drost@auckland.ac.nz
- gus.ellerm@auckland.ac.nz
- <https://github.com/LivePublication>