

Provenance Requirements for SKA

Jesús Salgado -SKA Regional Centre Architect Rosie Bolton - Head of Data Operations

SKA1-low – the SKA's low-frequency instrument

will revolutionise our understanding of the Universe. It will have a uniquely distributed character: one observatory operating two telescopes on three continents. Construction of the SKA will be phased and work is currently focused on the first phase named SKA1.



Compared to LOFAR Netherlands, the current

best similar instrument in the world

resolution sensitive

better SKA-low's resolution will

135x

the survey

speed

be similar to LOFAR.



SQUARE KILOMETRE ARRAY

SKA1-mid – the SKA's mid-frequency instrument

character: one observatory operating two telescopes on three continents. Construction of the SKA will be phased and work is currently focused on the first phase named SKA1. corresponding to a fraction of the full SKA. SKA1 will include two instruments – SKA1-mid





www.skatelescope.org 💆 @SKA_telescope 🥤 SKAtelescope 🔞 ska_telescope YouTube Square Kilometre Array in ska-organisation





*

Data Mesh: Domain Oriented



Domain Data as a product



- Discoverable
- Addressable
- Self-describing
- Trustworthy
- Secure
- InterOperable

Data Product in Domain





SKA Regional Centre Capabilities Blueprint







SRC Network global capabilities



Every node is an instance of the blueprint Interconnections are done using agreed APIs, using FAIR and VO protocols where available Collectively meet the needs of the global community of SKA users Anticipate heterogeneous SRCs, with different strengths

SRCNet principles: Use of Standards

- Build SKA science archive around FAIR and IVOA standards
- Ensure interoperability with other archives and other experiments
- Strong adherence to the FAIR principles
- Give credit appropriately to all contributors to a team







SRCNet principles: Collaboration and Reproducibility

- Most SKA projects will be collaborative
- SRCs will provide collaborative tools
 - Sharing components
 - Single Sign-on
- Support to workflows
- Provenance metadata
- Science Reproducibility at the level of workflows is essential as data should not be downloaded
- 43 requirements talking about provenance for the SRC Net, 26 L1 level



SRC Net Provenance Requirements Examples

- SRC-102 Interactive Analysis: Interactive analysis actions should be permit the serialization in the form of reproducible workflows
- SRC-169 Qualified references to other (meta)data within ADPs: Advanced Data Product would contain provenance metadata to characterize its generation
- SRC-138 and similar requirements. Provenance should be maintained for all products (including movements within the network)
- SRC-118 Uniform vocabulary in the archive: Provenance is part of the vocabulary shared by all the SRC Net node
- SRC-22 Reproducibility of results
- SRC-108 Long term preservation of metadata and provenance metadata for Observatory and Advanced Data Products

SRC Net Provenance non-written yet Requirements

- Provenance mapper
- Workflow editor and visualizer (e.g. Taberna)
- Support to needed execution framework
- Full characterization of software elements
 - Software Versioning
 - Software Access
 - Input/Output parameters/results
 - Error handling
 - Execution
 - Containers?
 - Specific Hardware and execution needs
 - Data Domain



Replicability vs Reproducibility

- **Reproducibility**: Ability to repeat an experiment with minor differences to the original experiment, but still achieving the same qualitative result.
 - E.g. Notebooks or pipelines
- **Replicability**: Ability to exactly reproduce results by running exactly the same experiment.
 - No changes in parameters or environment
 - Frozen data
- First one, quite important to generate pipelines, second to replicate results of a given experiment/paper
- Probably both things are needed for SKA with different level or priority

Mapping the reproducibility space



Provenance and data differencing for workflow reproducibility analysis Missier, P.; Woodman, S.; Hiden, H.; and Watson, P. Concurrency and Computation: Practice and Experience, . 2013

Reproducibility in Cloud

Listing 1.1. HyperFlow workflow description (fragment).

```
1 processes:
     "name": "alignment to reference",
2
     "function": "k8sCommand",
3
     "config": {
4
       "executor": {
5
         "executable": "bwa-wrapper",
6
         "args": [ "mem", "-t", "2", "-M",
7
           "Gmax 275 v2.0.fa",
8
           "USB-001 1. fastg",
9
           "USB-001 2.fastq"
10
11
         "cpuRequest": "1",
12
         "memRequest": "500M",
13
         "stdout": "20180321-083514-USB-001 aligned reads.sam"
14
15
16
     "ins": [1,2,3,4,5,6,7,8,9,10,11,12],
17
     "outs": [13]
18
   \{, \dots \},
19
   signals: [ {
20
     "name": "Gmax 275 v2.0.fa",
21
     "type": "file",
22
     "size": "990744229".
23
     "md5sum": "3aa6cf1962f5260cf1405e82efb25c71"
24
25
   }, ...
```

Reproducibility of Computational Experiments on Kubernetes-Managed Container Clouds with HyperFlow, V. Krzhizhanovskaya et al <u>Computational Science – ICCS 2020</u>. 2020 May 26; 12137: 220–233. Published online 2020 May 26. doi: <u>10.1007/978-3-030-50371-0_16</u>

Other thoughts

- How to combine data from other missions/data providers and maintain reproducibility?
 - Copy at SRC Net?
- Workflow Serialization Language (CWL?)
- Workflows simplification?
 - Graph analysis and proposal of more optimal execution
 - Analysis of resources
- How to combine provenance from different centres
 - Provenance discovery (Prov-TAP is enough?)
- Specific provenance fields for execution through Provenance extension?

Thanks for your attention

•

 \bullet

•

۲

۲