



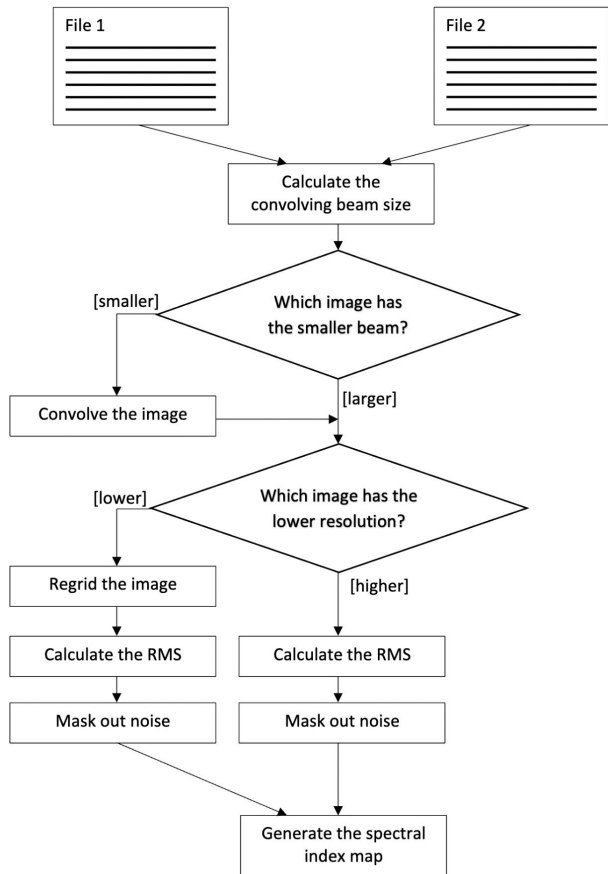
Tangerine team update (SKA Science Platform Prototyping)

Vision and Progress

Arpan Das(EPFL)
Rohit Sharma (FHNW)

Spectral Index Workflow

This workflow describes how to combine two images taken at different frequencies of the same source into a spectral index map.



- Brief overview of how existing workflow packages run and how the integration tools supplied by each workflow
- Not to give an in depth functional analysis of each workflow package
- This analysis is designed to give a more technical overview of how each workflow package could be used within the SRC Net

The packages studied are:

- Nextflow
- Snakemake
- Daliuge
- CWL
- WDL
- Ruffus
- Yadage
- Stimela

Nextflow:

NextFlow uses a scripting language similar to CWL, based upon a language called 'groovy' (an extension of java), and provides out-of-the-box executors for GridEngine, SLURM, LSF, PBS, Moab and HTCondor batch schedulers, and for Kubernetes, Amazon AWS, Google Cloud, and MS Azure platforms.

The NextFlow application offers the following commands:

- clean - clean up a project cache and work directories
- clone - clone a project into a folder
- config - print a project configuration
- console - launch NextFlow interactive console
- drop - delete the local copy of a project
- help - print the usage help for a command
- info - print project and system runtime information
- kuberun - execute a workflow in a Kubernetes cluster (experimental)
- list - list all downloaded projects
- log - print executions log and runtime info
- pull - download or update a project
- run - execute a pipeline project
- self-update - update NextFlow runtime to the latest available version
- view - view project script file(s)

Basic Concepts

- Designed around Linux Platform for data science applications
- Linux provide many simple but powerful CL and scripting tools, when chained together facilitate complex data manipulations
- Nextflow extends this approach, adding ability to define complex program and a high-level parallel computational environment based on dataflow programming model

Nextflow Pipeline Script are made by joining different processes

- Processes are executed independently and isolated from one another
- Each process can be written in any scriptable language which is executable on Linux platform (Bash, Perl, Ruby, Python etc..)
- They do share a common writable state
- The processes can communicate via asynchronous queues, called channels
- Any Process can define one or more channels as input and output

Features

Nextflow is built around the idea that Linux is the *lingua franca* of data science.

Fast prototyping

Nextflow allows you to write a computational pipeline by making it simpler to put together many different tasks.

You may reuse your existing scripts and tools and you don't need to learn a new language or API to start using it.

Portable

Nextflow provides an abstraction layer between your pipeline's logic and the execution layer, so that it can be executed on multiple platforms without it changing.

It provides out of the box executors for GridEngine, SLURM, LSF, PBS, Moab and HTCondor batch schedulers and for [Kubernetes](#), [Amazon AWS](#), [Google Cloud](#) and [Microsoft Azure](#) platforms.

Continuous checkpoints

All the intermediate results produced during the pipeline execution are automatically tracked.

This allows you to resume its execution, from the last successfully executed step, no matter what the reason was for it stopping.

Reproducibility

Nextflow supports [Docker](#) and [Singularity](#) containers technology.

This, along with the integration of the [GitHub](#) code sharing platform, allows you to write self-contained pipelines, manage versions and to rapidly reproduce any former configuration.

Unified parallelism

Nextflow is based on the *dataflow* programming model which greatly simplifies writing complex distributed pipelines.

Parallelisation is implicitly defined by the processes input and output declarations. The resulting applications are inherently parallel and can scale-up or scale-out, transparently, without having to adapt to a specific platform architecture.

Stream oriented

Nextflow extends the Unix pipes model with a fluent DSL, allowing you to handle complex stream interactions easily.

It promotes a programming approach, based on functional composition, that results in resilient and easily reproducible pipelines.

```

// Declare syntax version
nextflow.enable.dsl=2
// Script parameters
params.query = "/some/data/sample.fa"
params.db = "/some/path/pdb"

process processOne {
    input:
        path query
        path db
    output:
        path "pathOne.txt"

    // process is defined here.
    """
    """
}

process processTwo {
    input:
        path in_path
    output:
        path "pathTwo.txt"

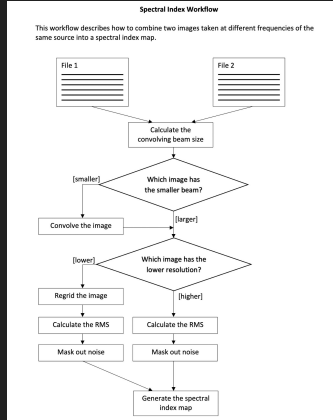
    // process is defined here.
    """
    """
}

workflow {
    def query_ch = Channel.fromPath( params.query )
    processOne( query_ch, params.db ) | processTwo | view
}

```

- The pipe between the two processes forwards the output from one process to the inputs of the following one
- Creates two processes (one for each function) and a channel (query_ch)
- Both processes will be started at the same time and they will listen to their respective input channels
- When processOne emit a value, processTwo will receive it
- Inside the triple quotes the user writes a bash script, which can kick off a script in python, R, or other languages
- To use something other than bash, the first line of the script should be `#!/usr/bin/perl`, or `#!/usr/bin/python`, etc
- To submit jobs via SLURM the process.executor in the nextflow.config file should be 'slurm'
- NextFlow integrates with github, and if a pipeline is not found locally by the executor it will automatically look for a public github repository with the same name, download it, and execute it.

Spectral Index Workflow



```

1 # MEASUREMENT_SET = file ~/measurement-sets/30391_cts.ksmic.spw.ms"
10 SCRIPT_IMAGE_0 = file ~/scripts/make_image_channel_0.py"
11 SCRIPT_IMAGE_03 = file ~/scripts/make_image_channel_03.py"
12 SCRIPT_BEAM_SIZE_0 = file ~/scripts/get_beam_size.py"
13 SCRIPT_BEAM_SIZE_03 = file ~/scripts/get_beam_size.py"
14 SCRIPT_SPECTRAL_INDEX = file ~/scripts/make_spectral_index_map.py"
15 SCRIPT_EXPORT_FITS = file ~/scripts/export_fits.py"
16 SCRIPT_CONVOLVE = file ~/scripts/convolve.py"
17
18 // -----
19 //
20 // processes
21 // -----
22 //
23 //
24 // make_image_channel_0
25 //
26 // Modified: 26/04/2023
27 //
28 // Make a casa image from channel 0.
29 //
30 //
31 // process make_image_channel_0
32 //
33 //
34 //
35 //
36 // container 'auiaphub/casa:6.5.2'
37 //
38 // input:
39 // val measurement_set
40 //
41 // NOTE: the exit keywords doesn't appear to work for type file, so I have to use type path.
42 //
43 // output:
44 // file 'x_image', emit: clean_image_0
45 // path '30391_cts_spx_0_ch_0.l_image', emit: clean_image
46 //
47 //
48 //
49 // docker run -v /home/tangerine/scripts/scripts -v /home/tangerine/measurement-sets/measurement-sets -v
50 //
51 //
52 //
53 // // make_image_channel_0
54 //
55 //
56 // make_image_channel_03
57 //
58 // Modified: 26/04/2023
59 //
60 //
61 // Make a casa image from channel 03.
62 //
63 //
  
```

```

137 process get_convolution_beam_size
138 {
139 //
140 // NOTE: for some reason I have to make these variables of type val instead of path or file, or otherwise
141 // does not pull in the files from the previous working directory.
142 //
143 // input:
144 // val beam_size_0_path
145 // val beam_size_1_path
146 //
147 // NOTE: the exit keywords doesn't appear to work for type file, so I have to use type path.
148 //
149 // output:
150 // path 'beam-size', emit: beam_size
151 // val beam_size
152 //
153 // NOTE: the exec block runs groovy script, the command task.workDir.resolve() doesn't work inside a sed
154 // strangely, inside an exec or script block the working directory is the directory from which the
155 // bash script (bounded by three double quotes) the working directory is the task directory of the
156 // task.workDir.resolve() to add the full path to the task directory.
157 //
158 // exec:
159 //
160 // get the beam sizes from the files.
161 // beam_size_0 = file task.workDir.resolve(beam_size_0_path).text.toFloat()
162 // beam_size_1 = file task.workDir.resolve(beam_size_1_path).text.toFloat()
163 //
164 // which one is the larger? update the return value.
165 // beam_size = beam_size_0
166 // if (beam_size_1 > beam_size_0)
167 // {
168 //     beam_size = beam_size_1
169 // }
170 //
171 // Increase the beam size by 2% so that we can do the convolution.
172 // beam_size = beam_size * 1.02
173 //
174 // NOTE: groovy has a max() function, but I couldn't get it to work.
175 // beam_size = max(beam_size_0, beam_size_1)
176 //
177 // write to an output file.
178 // NOTE: using a bash script that does 'echo ... >> beam-size' doesn't work for some reason, the
179 // task.workDir.resolve('beam-size').text = beam_size.toString()
180 //
181 // println "convolving to a circular beam of size " + beam_size.toString() + " arcsec"
  
```

```

271 process export_fits
272 {
273 //
274 // container 'auiaphub/casa:6.5.2'
275 //
276 // stageInMode 'copy'
277 //
278 // input:
279 // file casa_map
280 //
281 // output:
282 // file '30391_spectral_index.fits'
283 //
284 //
285 // docker run -v /home/tangerine/scripts/scripts -v /home/tangerine/works/work/58/7153760a74ed39d
286 //
287 //
288 // docker run -v /home/tangerine/scripts/scripts -v ./work auiaphub/casa:6.5.2 /casa-6.5.2-20-pys.6/bin.
289 //
290 //
291 // export_fits
292 //
293 // copy_required_images
294 //
295 // Modified: 26/04/2023
296 //
297 // Copy the casa and fits maps to the output folder (~images/).
298 //
299 //
300 // process copy_required_images
301 //
302 //
303 // input:
304 // file casa_map
305 // file fits_map
306 //
307 //
308 //
309 //
310 // copy_required_images
311 //
312 // -----
313 //
314 // workflow
315 // -----
316 //
317 //
318 // Modified: 26/04/2023
319 //
320 //
321 // Construct a spectral index map from channels 0 and 03, and export the casa image to fits format.
322 // Both the casa and fits maps are written to the ~/images/ directory.
323 //
324 //
  
```

```

325 workflow
326 {
327 //
328 // make images from channels 0 and 03 (in parallel)
329 // make_image_channel_0( MEASUREMENT_SET )
330 // make_image_channel_03( MEASUREMENT_SET )
331 //
332 // get the major and minor axes of the psfs.
333 //
334 // majorAxis = psfhead(image_name = 'image_image', mode = 'get', hdkey = 'maj')
335 // minorAxis = psfhead(image_name = 'image_image', mode = 'get', hdkey = 'min')
336 //
337 // the returned format is 'unit': 'arcsec', 'value': 18.5794627799961, so to get the value and unit I
338 //
339 // value = ps.getval(majorAxis)
340 // unit = ps.getunit(majorAxis)
341 //
342 // get the major axis restoring beam size for these images.
343 // get_beam_size_0( make_image_channel_0.out.clean_image )
344 // get_beam_size_03( make_image_channel_03.out.clean_image )
345 //
346 // get the larger of the two beam sizes.
347 // beam_size = get_convolution_beam_size( get_beam_size_0.out.beam_size, get_beam_size_03.out.beam_size )
348 //
349 // convolve image 0 with a Gaussian.
350 // convolve_0( make_image_channel_0.out.clean_image, beam_size )
351 //
352 // convolve image 1 with a Gaussian.
353 // convolve_03( make_image_channel_03.out.clean_image, beam_size )
354 //
355 // make a spectral index map from our two images.
356 // map_spectral_index = make_spectral_index_map( convolve_0.out.smooth_image, convolve_03.out.smooth_image )
357 // map_fits = export_fits( map_spectral_index )
358 //
359 // copy any files we need to the images directory.
360 // copy_required_images( map_spectral_index, map_fits )
361 //
362 // // workflow
  
```

Nextflow
Workflow

Stimela:

Stimela (the IsiZulu word for a train) is a platform independent radio interferometry scripting framework based on Python and the containerization technologies that now comes standard with all major Linux distributions.

Currently supports:

- Podman
 - Docker
 - Singularity
 - uDocker
-
- ❖ In this framework, radio interferometry related tasks such as imaging, calibration and data synthesis are executed in containers.
 - ❖ The packages that perform these tasks are Python modules.
 - ❖ Stimela does not do any data processing, synthesis or analysis but offers a unified Pythonic interface to packages that perform these tasks.

The primary aims of Stimela is to provide the following services to the Radio Astronomy community:

- A user friendly and modular scripting environment that gives general users easy access to novel radio interferometry calibration, imaging, and synthesis packages.
- Simplified installation and production deployment. All the software available to the Stimela user is prebuild and available on Docker Hub.

The production environment is fixed. The versions of the interfaced software products is fixed for a particular release version of Stimela. This means one can roll back and forward with ease and ensures that the production environment is always verbatim with that used for the original reduction work.

Stimela is centred around two sets of images:

- Base images, which have the required software tools installed in them. The base images can either be built locally (on the host machine) or pulled from Docker hub.
- Very light weight executor images (a.k.a cab images) based on the base images, these perform radio interferometry related tasks like imaging, data synthesis, and calibration. The executor images are built locally.

Base images:

- stimela/meqtrees - MeqTrees calibration/simulation tool
- stimela/lwimager - Uses the casarest based lwimager tool for imaging and deconvolution
- stimela/wsclean - WSClean imaging tool
- stimela/simms - Uses CASA simulate tool to create a simulated (empty) MS
- stimela/tigger - Tools for managing and manipulating analytic sky models (Gaussian and point sources)
- stimela/aoflagger - Automated RFI flagging tool
- stimela/casa - CASA
- stimela/lofar - Lofar
- stimela/sourcery - Source finding and source characterisation tool
- stimela/msutils - Convenience functions for manipulating MSs

Executors (a.k.a 'cabs'):

These images are generally pre-loaded with Python scripts that perform a specified task (e.g calibrating a visibility dataset). A stimela cab image takes some input as well a set of instructions, performs some task, then returns the output. The following are examples of available tasks:

- cab/simms
- cab/simulator
- cab/calibrator
- cab/lwimager
- cab/wsclean
- cab/tigger_convert
- cab/tigger_restore
- cab/tigger_tag
- cab/specfit
- cab/sourcery
- cab/autoflagger (AOFlagger)
- cab/flagms
- cab/casa_{clean, gaincal, bandpass, etc.}
- cab/ddfacet
- cab/cubical
- cab/tricolour (Tricolour)

SRC Science Analysis Platform Vision Document

de Boer, J.¹, Cimpan, I.², Das, A.³, Fabbro, S.⁴, Grange, Y. G.^{1*}, Hardcastle, M. J.⁵,
Sharma, R.⁶, Skipper, C. J.², Swinbank, J. D.¹, Webster, B.⁵

¹ASTRON, the Netherlands Institute for Radio Astronomy, Oude Hoogeveensedijk 4,7991 PD Dwingeloo,
The Netherlands

²Jodrell Bank Centre for Astrophysics, Alan Turing Building, The University of Manchester, Manchester,
M13 9PL, UK

³École polytechnique fédérale de Lausanne, Rte Cantonale, 1015 Lausanne, Switzerland

⁴NRC Herzberg Astronomy and Astrophysics, 5071 West Saanich Road, Victoria, BC V9E 2E7, Canada

⁵Centre for Astrophysics Research, University of Hertfordshire, College Lane, Hatfield, AL10 9AB, UK

⁶Fachhochschule Nordwestschweiz, Bahnhofstrasse 6, 5210 Windisch, Switzerland

*corresponding author

Abstract. This document describes the vision for the [Square Kilometer Array \(SKA\)](#) Regional Centres (SRC) Science Analysis Platform. It is intended to set the broad terms of reference for the platform and to provide guidance for both development teams and other stakeholders. Among the features and services that are expected to be included are data querying and discovery tools, some form of notebook interface, user-managed software environments, workflow management, and a comprehensive set of [application programming interfaces](#) (APIs) enabling access to all low-level platform functionality. This document is not a design specification, and the features and services described herein will be further refined, or could be discarded, at a later stage of development.

- Implementing external feedbacks in progress
- ~ 100 comments