# Bioinformatics Pipeline Frameworks

Stephen Kelly

NYU Langone Medical Center

November 13, 2017

# Advantages

- Standardized format

- Modular

- Consistency

- Automate the Easy Stuff... More Easily!

# snsxt

- https://github.com/NYU-Molecular-Pathology/snsxt
- http://snsxt.readthedocs.io/en/latest/snsxt.html

- Modular, object oriented, written in Python

- Easy wrapper around misc. analysis tasks

  - run Python code directly

  - run external scripts

- Easy to wrap entire pipelines

  - built-in wrapper around Igor's sns pipeline (original purpose)

- **Tailored for easy usage on NYULMC phoenix HPC cluster**

  - modules: Python 2.7+, pandoc 1.13, batteries included; **clone & run**

  - `$ snsxt/run.py –d analysis_dir/ –f fastq_dir/`

# Example sns task – sns wes pipeline

import base class

run external `sns` command

capture qsub jobs

```python
#!/usr/bin/env python
# -*- coding: utf-8 -*-

import os
from task_classes import SnsTask
from time import sleep

class SnsWes(SnsTask):
    """
    Run the sns wes analysis pipeline for unpaired variant calling on exome data
    """
    def __init__(self, analysis_dir, taskname = 'SnsWes', extra_handlers = None,
    **kwargs):
        """
        """
        SnsTask.__init__(self, analysis_dir = analysis_dir, taskname = taskname,
        extra_handlers = extra_handlers)

    def run(self, *args, **kwargs):
        """
        """
        expected_log_dir = os.path.join(self.analysis_dir, "logs-qsub")
        command = 'sns/run wes'
        run_cmd = self.run_sns_command(command = command)
        jobs = self.catch_sns_jobs(proc_stdout = run_cmd.proc_stdout, log_dir =
        expected_log_dir)

        # wait a few seconds to allow time for jobs to initialize
        sleep(10)
        return(jobs)
```

# External config



```
StartSns.yml                    ×
1   # ~~~~ REQUIRED TASK ITEMS ~~~~ #
2   # every sns_task should have these items
3
4   # name of the parent Python module
5   task_name: StartSns
6
7
8   # name of the sns output subdirectory from which to take input files
9   # input_dir: '.'
10  # ^ this will be ignored
11
12  # filename pattern to use for input file
13  # input_pattern: '*.dd.ra.rc.bam'
14
15  # or exact suffix to append to sample ID for input file
16  # input_suffix: ''
17
18  # name of the parent directory to use for the program output
19  # output_dir_name: Demo-QsubSampleTask
20  # i.e. analysis_dir/QC-Coverage-Custom will be used
21
22  # files in the `report_dir` associated with this sns_task; should end in '_report.Rmd'
23  report_files:
24
25
26  # ~~~~ ANALYSIS TASK ITEMS ~~~~ #
27  # use these if the task will operate on the analysis as a whole
28
29  # input_files:
30  # - 'baz.txt'
31
32  # files that should be output by the analysis task
33  output_files:
34    - settings.txt
35    - summary-combined.wes.csv
36    - samples.fastq-raw.csv
37
38  # files that should be sent in email output for the task
39  email_files:
40    - settings.txt
41    - summary-combined.wes.csv
42    - samples.fastq-raw.csv
43
44  # ~~~~ TASK SPECIFIC CUSTOM ITEMS ~~~~ #
45
```

**files to use for reporting** → (line 23)

**expected input files** → (line 29)

**expected output files** → (line 33)

**files to send in email output** → (line 39)
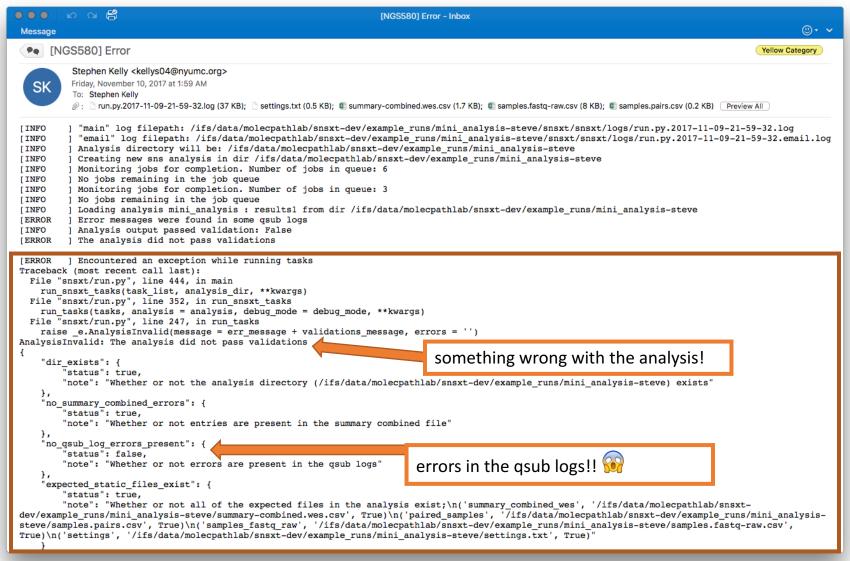
# YAML Task List



```yaml
default.yml                    ×
1   # this task list will perform a full, new
    analysis, from start to finish
2
3   # sns analysis setup tasks
4   sns:
5     StartSns:
6     SnsWes:
7     SnsWesPairsSnv:
8
9   # downstream analysis tasks
10  tasks:
11    Delly2:
12      qsub_wait: False
13    GATKDepthOfCoverageCustom:
14      qsub_wait: True
15    SummaryAvgCoverage:
16    HapMapVariantRef:
17
18  # compile the report for the analysis
19  setup_report: True
```

start a new `sns` analysis

downstream tasks

task name = Python class name

extra args for the task's `run()` method

compile report for the analysis

# Email output – oops it broke!

# Email output – yay it worked!

# Reporting

- Modular, extensible report framework
  - R Markdown + pandoc
- custom report per pipeline task
- parent report imports & compiles all child report docs

# Other Frameworks



- Advantages:
  - might have more features

- Disadvantages:
  - how to implement?
  - how to use?
  - how to apply to our Bioinformatic pipelines?

# Conclusion

- snsxt's advantages

  - already works on NYULMC phoenix HPC out of the box

  - no installation needed, single module adjustment

    - `module load python/2.7`

  - standardized format to wrap existing pipelines & analysis output

  - builds off of our current code base

  - extensible for new pipelines