# pyrpipe

*Release 0.0.1*

**Jun 04, 2021**

# Contents

**Author** Urminder Singh

**Date** Jun 04, 2021

**Version** 0.0.5

Introduction

pyrpipe (pronounced as "pyre-pipe") is a python package to easily develop computational pipelines in pure python, in an object-oriented manner. pyrpipe provides an easy-to-use object oriented framework to import any UNIX executable command or tool in python. pyrpipe also provides flexible and easy handling of tool options and parameters and can automatically load the tool parameters from .yaml files. This framework minimizes the commands user has to write, rather the tools are available as objects and are fully re-usable. All commands executed via pyrpipe are automatically logged extensively and could be compiled into reports using the pyrpipe_diagnostic tool.

To enable easy and easy processing of RNA-Seq data, we have implemented specialized classes build on top of the pyrpipe framework. These classes provide high level APIs to many popular RNA-Seq tools for easier and faster development of RNA-Seq pipelines. These pipelines can be fully customizable and users can easily add/replace the tools using the pyrpipe framework. Finally, pyrpipe can be used on local computers or on HPC environments and pyrpipe scripts can be easily integrated into workflow management systems such as Snakemake and Nextflow.

Get started with the tutorial: *Tutorial* or look at real data examples.

## 1.1 Key Features

- Import any UNIX executable command/tool in python

- Dry-runnable pipelines to check dependencies and commands before execution

- Flexible and robust handling of tool arguments and parameters

- Auto load parameters from .yaml files

- Easily override threads and memory options using global values

- Extensive logging and reports with MultiQC reports for bioinformatic pipelines

- Specify GNU make like targets and verify the integrity of the targets

- Automatically resume pipelines/jobs from where interrupted

- Easily integrated into workflow managers like Snakemake and NextFlow

## 1.2 Installation

To install the latest stable release via conda:

```
conda install -c bioconda pyrpipe
```

To install the latest stable release via pip:

```
pip install pyrpipe --upgrade
```

Install latest development version :

```
git clone https://github.com/urmi-21/pyrpipe.git
pip install -r pyrpipe/requirements.txt
pip install -e pyrpipe
```

See the *Installation notes* for details.

Examples and Case-Studies

Example usage and case-studies with real data is provided at GitHub.

## 2.1 Installation

**Note:** See *Create a new conda environment* in tutorial, to learn how to install pyrpipe and required tools in an conda environment.

Before installing pyrpipe, make sure conda channels are added in right order

```
conda config --add channels defaults
conda config --add channels bioconda
conda config --add channels conda-forge
```

pyrpipe is available through bioconda could be installed via:

```
conda install -c bioconda pyrpipe
```

pyrpipe is available on PyPI and could be installed via pip:

```
pip install pyrpipe
```

To install from source, clone the git repo:

```
git clone https://github.com/urmi-21/pyrpipe
```

Then install using pip:

```
pip install -r pyrpipe/requirements.txt
pip install -e pyrpipe
```

To run tests using pytest, from the pyrpipe root perform:

```
py.test
#or
py.test tests/<specific test file>
```

## 2.2 Tutorial

This tutorial provides an introduction to pyrpipe. New users can start here to get an idea of how to use pyrpipe APIs for RNA-Seq or other analysis. This tutorial is divided into smaller sections.

Section 1 describes the recommended way of installing pyrpipe and setting up conda environments.

Section 2 is an introduction to basic RNA-Seq API is provided via a simple RNA-Seq processing pipeline example.

Section 3 provides information on how to extend pipelines using third-party tools and integrate pyrpipe into snakemake for more scalable workflows.

Subsequent Sections detail pyrpipe_engine, pyrpipe_utils and pyrpipe_diagnostic, each of which provide helpful functionality.

### 2.2.1 Setting up the environment

First step is to setup the environment in a way such that reproducibility is ensured or maximized. In this tutorial we will use the conda environment manager to install python, pyrpipe and the required tools and dependencies into a single environment. **Note:** Conda must be installed on the system. For help with setting up conda, please see miniconda.

#### Create a new conda environment

To create a new conda environment with python 3.8 execute the following commands. We recommend sharing conda environment files with pipeline scripts to allow for reproducible analysis.

```
conda create -n pyrpipe python=3.8
```

Activate the newly created conda environment and install required tools

```
conda activate pyrpipe
conda install -c bioconda pyrpipe star=2.7.7a sra-tools=2.10.9 stringtie=2.1.4 trim-
→galore=0.6.6 orfipy=0.0.3 salmon=1.4.0
```

If the above command fails, please try adding conda channels (see commands below) in the right order and then try again.

```
conda config --add channels defaults
conda config --add channels bioconda
conda config --add channels conda-forge
```

#### Using conda environment in yaml files

We have provided a yaml file containing the conda packages required to reproduce pyrpipe environment. Users can also use this file to create a conda environment and run pyrpipe. To create a conda environment, use the pyrpipe_environment.yaml:

```
conda env create -f pyrpipe_environment.yml
```

Users can easily export and share their own conda environment yaml files containing information about the conda environment. To export any conda environment as yaml, run the following command

```
conda env export | grep -v "^prefix: " > environment.yml
```

To recreate the conda environment in the *environment.yml*, use

```
conda env create -f environment.yml
```

### Automated installation of required tools

We have also provided a utility to install required RNA-Seq tools via a single command:

```
pyrpipe_diagnostic build-tools
```

**Note:** Users must verify the versions of the tools installed in the conda environment.

### Setting up NCBI SRA-Tools

After installing sra-tools, please configure prefetch to save the downloads the the **public user-repository**. This will ensure that the prefetch command will download the data to the user defined directory. To do this

- Type *vdb-config -i* command in terminal to open the NCBI SRA-Tools configuration editor.
- Under the TOOLS tab, set prefetch downloads option to **public user-repository**

Users can easily test if SRA-Tools has been setup properly by invoking the following command

```
pyrpipe_diagnostic test
```

## 2.2.2 Basic RNA-Seq processing

After setting up the environment, one can import pyrpipe modules in python and start using it. In this example we will use the *A. thaliana* Paired-end RNA-Seq run *SRR976159*.

Many other examples are available on github

### Required files

We need to download the reference genome and annotation for *A. thalinia*. This can be done inside python script too. For simplicity we just download these using the *wget* command from the terminal.

Code for *A. thalinia* transcript assembly case study is available on github

```
1  wget ftp://ftp.ensemblgenomes.org/pub/release-46/plants/fasta/arabidopsis_thaliana/
   ↪dna/Arabidopsis_thaliana.TAIR10.dna.toplevel.fa.gz
2  gunzip Arabidopsis_thaliana.TAIR10.dna.toplevel.fa.gz
3  wget ftp://ftp.ensemblgenomes.org/pub/release-46/plants/gtf/arabidopsis_thaliana/
   ↪Arabidopsis_thaliana.TAIR10.46.gtf.gz
4  gunzip Arabidopsis_thaliana.TAIR10.46.gtf.gz
```

**Simple pipeline**

RNA-Seq processing is as easy as creating required objects and executing required functions. The following python script provides a basic example of using pyrpipe on publicly available RNA-Seq data.

```python
from pyrpipe import sra,qc,mapping,assembly
#define some vaiables
run_id='SRR976159'
working_dir='example_output'
gen='Arabidopsis_thaliana.TAIR10.dna.toplevel.fa'
ann='Arabidopsis_thaliana.TAIR10.46.gtf'
star_index='star_index/athaliana'
#initialize objects
#creates a star object to use with threads
star=mapping.Star(index=star_index,genome=gen,threads=4)
#use trim_galore for trimming
trim_galore=qc.Trimgalore()
#Stringtie for assembly
stringtie=assembly.Stringtie(guide=ann)
#create SRA object; this will download fastq if doesnt exist
srr_object=sra.SRA(run_id,directory=working_dir)
#create a pipeline using the objects
srr_object.trim(trim_galore).align(star).assemble(stringtie)

#The assembled transcripts are in srr_object.gtf
print('Final result',srr_object.gtf)
```

The above code defines a simple pipeline (in a single line: Line 18) that:

- Downloads fastq files from NCBI-SRA

- Uses Trim Galore for trimming

- Uses Star for alignemnt to ref genome

- Uses Stringtie for assembly

**A line by line explanation:**

1. Imports the required pyrpipe modules

3. Lines 3 to 7 defines the variables for reference files, *output directory*, and star index. The *output directory* will be used to store the downloaded RNA-Seq data and will be the default directory for all the results.

10. Creates a Star object. It takes index and genome as parameters. It will automatically verify the index and if an index is not found, it will use the genome to build one and save it to the index path provided.

12. Creates a Trimgalore object

14. Creates a Stringtie object

16. Creates an SRA object. This represents the RNA-Seq data. If the raw data is not available on disk it auto-downloads it via fasterq-dump.

18. This is the pipeline which describes a series of operations. The SRA class implements the *trim()*, *align()* and *assemble()* methods.

- *trim()* takes a qc type object, performs trimming via *qc.perform_qc* method and the trimmed fastq are updated in the SRA object

- *align()* takes a mappping type object, performs alignemnt via *mapping.perform_alignemnt* method. The result-ing bam file is stored in *SRA.bam_path*

- *assemble()* takes a assembly type object, performs assembly via *mapping.perform_assembly* method. The resulting gtf file is stored in *SRA.gtf*

### Executing the pipeline

To execute the pipeline defined above, save the python code in a file *pipeline.py*. The code can be executed as any python script using the python command:

```
python pipeline.py
```

Or it can be executed using the pyrpipe command and specifying the input script with the –in option

```
pyrpipe --in pipeline.py
```

One can specify pyrpipe specific options too

```
python pipeline.py --threads 10 --dry-run
    #OR
pyrpipe --in pipeline.py --threads 10 --dry-run
```

The above two commands are equivalent and specifies pyrpipe to use 10 threads. Thus 10 threads will be used for each of the tool except for STAR where we explicitly specied to use 4 threads during object creation.

The other option provided here is the *–dry-run* option and this option *turns off* the pyrpipe_engine and any command passed to the pyrpipe_engine is not actually executed but just displayed on screen and logged. During dry run the Runnable class also verifies the file dependencies (if any). More details are provided in the later chapters of the tutorial.

We recommend using the dry-run option before actually starting a job to make sure all parameters/dependencies are correct.

### Specifying tool parameters

pyrpipe supports auto-loading of tool parameters specified in .yaml files. The .yaml files must be stored is a directory and can be specified using the *–param-dir* option. The default value is *./params*. The files must be named as *<tool>.yaml*, for example star.yaml for STAR parameters. These parameters are loaded during object creation and user can easily override these during execution.

Create a directory *params* in the current directory and make a file *star.yaml* inside *params*. Add the following to *star.yaml* and rerun *pipeline.py* using the dry run option.

```
--outSAMtype: BAM Unsorted SortedByCoordinate
--outSAMunmapped: Within
--genomeLoad: NoSharedMemory
--chimSegmentMin: 15
--outSAMattributes: NH HI AS nM NM ch
--outSAMattrRGline: ID:rg1 SM:sm1
```

If you did everything correctly, you wil notice that now the STAR commands contain these specified parameters.

### Updating parameters dynamically

The parameters specified in the yaml will be replaced by any parameters provided during object creation. For example, consider the star.yaml specifying –runThreadN as 20

```
1  --outSAMtype: BAM Unsorted SortedByCoordinate
2  --outSAMunmapped: Within
3  --genomeLoad: NoSharedMemory
4  --runThreadN: 20
```

Now, consider creating a star object in the following scenarios

```
1  star1=Star(index='index') #will use 20 threads as mentioned in the yaml file
2  star2=Star(index='index',threads=5) #will use 5 threads
3  star3=Star(index='index',**{'--runThreadN':'10'}) #will use 10 threads
4  star4=Star(index='index') #initialized with 20 threads
5  star4.run(...,**{'--runThreadN':'10'}) #will use 10 threads for this particular run;
   →'--runThreadN':'20' remains in the star4 object
```

### 2.2.3 Importing tools into python

pyrpipe's *Runnable* can be used to import any Unix command into python. The *Runnable* class implements the *run()* method which checks required dependencies, monitors execution, execute commands via pyrpipe_engine, and verify target files. We will first download the *E. coli* genome file to use as input to orfipy.

```
1  wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/005/845/GCF_000005845.2_ASM584v2/
   →GCF_000005845.2_ASM584v2_genomic.fna.gz
2  gunzip GCF_000005845.2_ASM584v2_genomic.fna.gz
```

#### Basic example

In this tutorial we will consider the tool *orfipy*, a tool for fast and flexible ORFs, and import it in python.

```
1  from pyrpipe.runnable import Runnable
2  infile='GCF_000005845.2_ASM584v2_genomic.fna'
3  #create a Runnable object
4  orfipy=Runnable(command='orfipy')
5  #specify orfipy options; these can be specified into orfipy.yaml too
6  param={'--outdir':'orfipy_out','--procs':'3','--dna':'orfs.fa'}
7  orfipy.run(infile,**param)
```

Save the above script in *example.py* and execute it using *python example.py –dry-run*. pyrpipe should generate the orfipy command and display on screen during dry-run. Running it without *–dry-run* flag will generate the *orfipy_out/orfs.fa* file.

#### Requirements and targets

One can specify required dependencies and expected target files in the run() method Replacing the call to *run()* with the following will verify the required files and the target files. If command is interrupted, pyrpipe will scan for *Locked* taget files and resume from where the pipeline was interrupted.

```
1  orfipy.run(infile,requires=infile,target='orfipy_out/orfs.fa',**param)
```

### 2.2.4 Building APIs

Users can extend the Runnable class and specify classes dedicated to tools. Extra functionalities can be added by defining custom behaviour of classes. The RNA-Seq API built using pyrpipe the framework.

A small example is presented here to build a class for orfipy tool

```python
from pyrpipe import Runnable

from pyrpipe.runnable import Runnable
from pyrpipe import pyrpipe_utils as pu
from pyrpipe import sra
from pyrpipe import _threads,_mem,_force
import os


class Orfipy(Runnable):
    """
    Extends Runnable class
    Attributes
    ----------


    """
    def __init__(self,*args,threads=None,mem=None,**kwargs):
        """
        init an Orfipy object

        Parameters
        ----------
        *args : tuple
            Positional arguements to orfipy
        threads : int, optional
            Threads to use for orfipy. This will override the global --threads
    →parameter supplied to pyrpipe. The default is None.
        mem : int, optional
            Maximum memory to use in MB. The default is None.
        **kwargs : dict
            options for orfipy
        Returns
        -------
        None.


        """
        super().__init__(*args,command='orfipy',**kwargs)
        self._deps=[self._command]
        self._param_yaml='orfipy.yaml'
        #valid arguments for orfipy
        self._valid_args=['--min','--between-stops','--include-stop','--dna','--pep',
    →'--bed','--bed12','--procs','--chunk-size','--outdir'] #valid arguments for orfipy

        #resolve threads to use
        """
        orfipy parameter for threads is --procs
        if threads is passed in __init__() it will be used
        else if --procs is found in orfipy.yaml that will be used
        else if --procs is found in the passed **kwargs in __init__() it will be used
        else the default value i.e. _threads will be used
        if default value is None nothing will be done
        after the function, --procs and its value will be stored in self._kwargs, and
    →_threads variable will be stored in the Orfipy object.
        """
        self.resolve_parameter("--procs",threads,_threads,'_threads')
        #resolve memory to use
        """
```

(continues on next page)

```
54        default value is None--> if mem is not supplied don't make the self._mem␣
    ↪variable
55        """
56        self.resolve_parameter("--chunk-size",mem,None,'_mem')
57
58    ##now we write a custom function that can be used with an SRA object
59    def find_orfs(self,sra_object):
60        out_dir=sra_object.directory
61
62        out_file=os.path.join(out_dir,sra_object.srr_accession+"_ORFs.bed")
63
64        if not _force and pu.check_files_exist(out_file):
65            pu.print_green('Target files {} already exist.'.format(out_file))
66            return out_file
67
68        #In this example use orfipy on only first fastq file
69        internal_args=(sra_object.fastq_path,)
70        internal_kwargs={"--bed":sra_object.srr_accession+"_ORFs.bed","--outdir":out_
    ↪dir}
71
72        #call run
73        status=self.run(*internal_args,objectid=sra_object.srr_accession,target=out_
    ↪file,**internal_kwargs)
74
75        if status:
76            return out_file
77
78        return ""
```

The above class, Orfipy, we created can be used directly with SRA type objects via the find_orfs function. We can also still use the *run()* method and provide any input to orfipy.

```
1   #create object
2   orfipy=Orfipy()
3   #use run()
4   orfipy.run('test.fa',**{'--dna':'d.fa','--outdir':'of_out'},requires='test.fa',target=
    ↪'of_out/d.fa')
5
6   #use the api function to work with SRA
7   srr=sra.SRA('SRR9257212')
8   orfipy.find_orfs(srr)
```

Now try passing orfipy parameters from orfipy.yaml file. Create *params/orfipy.yaml* and add the following options into it.

```
1   --min: 36
2   --between-stops: True
3   --include-stop: True
```

Now re-run the python code and it will automatically read orfipy options from *./params/orfipy.yaml*.

## 2.2.5 Snakemake example

Since Snakemake direclty supports python, pyrpipe libraries can be driectly imorted into snakemake. Advantage of using a workflow manager like Snakemake is that it can handle parallel job-scheduling and scale jobs on clusters.

---

## Basic RNA-Seq example

A basic example of directly using pyrpipe with Snakemake is provided here. This example uses yeast RNA-Seq samples from the GEO accession GSE132425.

First run the following bash script to download the yeast reference genome from Ensembl. The last command also generated a star index with the downloaded genome under the *refdatayeast_index* directory.

```bash
#!/bin/bash
mkdir -p refdata
wget ftp://ftp.ensemblgenomes.org/pub/fungi/release-49/fasta/saccharomyces_cerevisiae/
→dna/Saccharomyces_cerevisiae.R64-1-1.dna.toplevel.fa.gz -O refdata/Saccharomyces_
→cerevisiae.R64-1-1.dna.toplevel.fa.gz
wget ftp://ftp.ensemblgenomes.org/pub/fungi/release-49/gff3/saccharomyces_cerevisiae/
→Saccharomyces_cerevisiae.R64-1-1.49.gff3.gz -O refdata/Saccharomyces_cerevisiae.R64-
→1-1.49.gff3.gz
cd refdata
gunzip -f *.gz
#run star index
mkdir -p yeast_index
STAR --runThreadN 4 --runMode genomeGenerate --genomeDir yeast_index --
→genomeFastaFiles Saccharomyces_cerevisiae.R64-1-1.dna.toplevel.fa --
→genomeSAindexNbases 10
```

Next, we create the following *snakefile*.

```python
import yaml
import sys
import os
from pyrpipe import sra,qc,mapping,assembly

####Read config#####
configfile: "config.yaml"
DIR = config['DIR']
THREADS=config['THREADS']
##check required files
GENOME= config['genome']
GTF=config['gtf']
#####Read SRR ids######
with open ("runids.txt") as f:
    SRR=f.read().splitlines()

#Create pyrpipe objects
#parameters defined in ./params will be automatically loaded, threads will be␣
→replaced with the supplied value
tg=qc.Trimgalore(threads=THREADS)
star=mapping.Star(threads=THREADS)
st=assembly.Stringtie(threads=THREADS)

rule all:
    input:
            expand("{wd}/{sample}/Aligned.sortedByCoord.out_star_stringtie.gtf",
→sample=SRR,wd=DIR),

rule process:
    output:
            gtf="{wd}/{sample}/Aligned.sortedByCoord.out_star_stringtie.gtf"
    run:
```

(continues on next page)

```
31              gtffile=str(output.gtf)
32              srrid=gtffile.split("/")[1]
33              sra.SRA(srrid,directory=DIR).trim(tg).align(star).assemble(st)
```

The above snakefile requires some additional files. A config.yaml file contains path to data and threads information.

```
DIR: "results"
THREADS: 5
genome: "refdata/Saccharomyces_cerevisiae.R64-1-1.dna.toplevel.fa"
gtf: "refdata/Saccharomyces_cerevisiae.R64-1-1.49.gff3"
```

Next the snakefile reads a file, *runids.txt* containing SRR accessions.

```
SRR9257163
SRR9257164
SRR9257165
```

Finally, we need to provide tool parameters for pyrpipe. Create a file *./params/star.yaml* and specify the index in it

```
--genomeDir: ./refdata/yeast_index/
```

Now the snakefile could be run using the snakemake command e.g *snakemake -j 8*

### pyrpipe_conf.yaml

Users can create a yaml file, *pyrpipe_conf.yaml*, to specify pyrpipe parameters, instead of directly passing them as command line arguments. When the *pyrpipe_conf.yaml* is found the pyrpipe specific arguments passed via the command-line are ignored. An example of *pyrpipe_conf.yaml* is shown below with the pyrpipe default values

```
dry: False          # Only print pyrpipe's commands and not execute anything through
→pyrpipe_engine module
threads: None       # Set the number of threads to use
force: False        # Force execution of commands if their target files already exist
params_dir: ./params # Directory containing parameters
logs: True          # Enable or disable pyrpipe logs
logs_dir: ./pyrpipe_logs    # Directory to save logs
verbose: False      # Display pyrpipe messages
memory: None        # Set memory to use (in GB)
safe: False         # Disable file deletion via pyrpipe commands
multiqc: False      # Automatically run multiqc after analysis
```

## 2.2.6 RNA-Seq API

pyrpipe implements specialized classes for RNA-Seq processing. These classes are defined into different modules, each designed to capture steps integral to analysis of RNA-Seq data –from downloading raw data to trimming, alignment and assembly or quantification. Each of these module implements classes coressponding to a RNA-Seq tool. These classes extend the *Runnable* class. Specialized functions are implemented such that analysis of RNA-Seq data is intuitive and easy to code. The following table provide details about the pyrpipe's RNA-Seq related modules.

| Module | Class | Purpose |
|--------|-------|---------|
| assembly | Assembly | Abstract class to represent Assebler type |
| assembly | Stringtie | API to Stringtie |
| assembly | Cufflinks | API to Cufflinks |
| mapping | Aligner | Abstract class for Aligner type |
| mapping | Star | API to Star |
| mapping | Bowtie2 | API to Bowtie2 |
| mapping | Hisat2 | API to Hisat2 |
| qc | RNASeqQC | Abstract class for RNASeqQC type (quality control and trimming) |
| qc | Trimgalore | API to Trim Galore |
| qc | BBmap | API to bbduk.sh |
| quant | Quant | Abstract class for Quantification type |
| quant | Salmon | API to Salmon |
| quant | Kallisto | API to Kallisto |
| sra | SRA | Class to represent RNA-Seq data and API to NCBI SRA-Tools |
| tools | Samtools | API to Samtools and other commonly used tools |

## The SRA class

The SRA class contained in the sra module captures RNA-Seq data. It can automatically download RNA-Seq data from the NCBI SRA servers via the prefetch command. The SRA constructor can take the SRR accession, path to fastq or sra file as arguments.

The main attributes and functions are defined the following table.

| Attribute | Description |
|-----------|-------------|
| fastq_path | Path to the fastq file. If single end this is the only fastq file. |
| fastq2_path | Path to the second fastq file for paired end data. |
| sra_path | Path to the sra file |
| srr_accession | The SRR accession for RNA-Seq run |
| layout | RNA-Seq layout, auto determined by SRA class. |
| bam_path | Path to bam file after running the align() function |
| gtf | Path to the gtf file after running assemble() function |

| Function | Description |
|---|---|
| __init__() | This is the constructor. It can take SRR accession, path to fastq files, or sra file as input. If accession if provided as input the files are downloaded via prefetch if they aren't preset on disk. It will automatically handle single-end and paired-end data. |
| download_sra() | This function downloads the sra file via prefetch. |
| download_fastq() | This function runs fasterq-dump on the sra file downloaded via prefetch. |
| sra_exists() | Check if fastq sra files are present |
| fastq_exists() | Check if fastq file exists |
| delete_sra() | Delete the sra file |
| delete_fastq() | Delete the fastq files |
| trim() | This function takes a RNASeqQC type object and performs trimming. The trimmed fastq files are then stored in fastq_path and fastq2_path. |
| align() | This function takes an Aligner type object and performs read alignemnt. The BAM file returned is stored in bam_path attribute. |
| assemble() | This function takes an Assembly type object and performs transcript assembly. The result is soted on the SRA object as gtf attributes |
| quant() | This function takes a Quantification type object and performs quant. |

### The RNASeqQC class

The RNASeqQC is an abstract class defined in the qc module. RNASeqQC class extends the Runnable class and thus has all the attributes as in the Runnable class. Classes Trimgalore and BBmap extends RNASeqQC class and share following attributes and functions.

| Attribute | Description |
|---|---|
| _category | Represents the type: "RNASeqQC" |

| Function | Description |
|---|---|
| __init__() | The constructor function |
| perform_qc() | Takes a SRA object, performs qc and returns path to resultant fastq files |

### The Aligner class

The Aligner is an abstract class defined in the mapping module. Aligner class extends the Runnable class and thus has all the attributes as in the Runnable class. Classes Star, Hisat2 and Bowtie2 extends the Aligner class and share following attributes and functions.

| Attribute | Description |
|---|---|
| _category | Represents the type: "Aligner" |
| index | Index used by the aligner tool |
| genome | Reference genome used by the tool |

| Function | Description |
|---|---|
| __init__() | The constructor function |
| build_index() | Build an index for the aligner tool using the *genome*. |
| check_index() | Checks if the index is valid |
| perform_alignment() | Takes a sra object, performs alignemnt and returns path to the bam file |

### The Assembly class

The Assembly is an abstract class defined in the assembly module. Assembly class extends the Runnable class and thus has all the attributes as in the Runnable class. Classes Stringtie and Cufflinks extends the Assembly class and share following attributes and functions.

| Attribute | Description |
|---|---|
| _category | Represents the type: "Assembler" |

| Function | Description |
|---|---|
| __init__() | The constructor function |
| perform_assembly() | Takes a SRA object, performs transcript assembly and returns path to resultant gtf/gff files |

### The Quant class

The Quant is an abstract class defined in the quant module. Quant class extends the Runnable class and thus has all the attributes as in the Runnable class. Classes Salmon and Kallisto extends the Quant class and share following attributes and functions.

| Attribute | Description |
|---|---|
| _category | Represents the type: "Quantification" |
| index | Index used by the aligner tool |
| transcriptome | Reference transcriptome used by the tool |

| Function | Description |
|---|---|
| __init__() | The constructor function |
| build_index() | Build an index for the quantification tool using the *transcriptome*. |
| check_index() | Checks if the index is valid |
| perform_quant() | Takes a sra object, performs quantification and returns path to the quantification results file |

### 2.2.7 pyrpipe_engine module

The pyrpipe_engine module is the key module responsible for handling execution of all the Unix commands. The Runnable class calls the execute_comand() function in the pyrpipe_engine module. All the commands executed via the pyrpipe_engine module are automatically logged. All the functions responsible for executing Unix command are "decorated" with the dryable method, which allows using the –dry-run flag on any function using the pyrpipe_engine module.

The pyrpipe_engine can be directly used to execute Unix command or to import output of a Unix command in python. The functions defined in the pyrpipe_engine module are described below.

| Function | Description |
|---|---|
| dryable() | A decorater to make Unix functions dry when –dry-run is specified |
| parse_cmd() | Parse a Unix command and return command as string |
| get_shell_output() | Execute a command and return return code and output. These commands are not logged |
| get_return_status() | Execute command and return status of the command |
| execute_commandRealtime() | Execute a command and print output in realtime |
| execute_command() | Execute a command and log the status. This is the function used by the Runnable class. |
| is_paired() | Check is sra file paired or single end. Uses fastq-dump |
| get_program_path() | Return path to a Unix command |
| check_dependencies() | Check a list of dependencies are present in Unix path |
| delete_files() | Delete files, rm command |
| move_file() | Moves files, mv command |

## 2.2.8 pyrpipe_utils module

The pyrpipe_utils module define multiple helpful functions. Functions defined in the pyrpipe_utils modules are extensively used throughout pyrpipe modules. Users can directly utilize these fuctions in their code and expedite development. A description of these functions is provided below.

| Function | Description |
|---|---|
| pyrpipe_print() | Prints in color |
| get_timestamp() | Return current timestamp |
| check_paths_exist() | Return true is paths are valid |
| check_files_exist() | Return True if files are valid |
| check_hisatindex() | Verify valid Hisat2 index |
| check_kallistoindex() | Verify kallisto index |
| check_salmonindex() | Verify salmon index |
| check_starindex() | Verify STAR index |
| check_bowtie2index() | Verify Bowtie2 index |
| get_file_size() | Return file size in human readable format |
| parse_java_args() | Parse tool options in JAVA style format |
| parse_unix_args() | Parse tool options in Unix style format |
| get_file_directory() | Return a file's directory |
| get_filename() | Return filename with extension |
| get_fileext() | Return file extension |
| get_file_basename() | Return filename, without extension |
| mkdir() | Create a directory |
| get_union() | Return union of lists |
| find_files() | Search files using regex patterns |
| get_mdf() | Compute and return MD5 checksum of a file |

## 2.2.9 pyrpipe_diagnostic

The pyrpipe_diagnostic command allows users to easily examine pyrpipe logs. It provides several options.

1. report: This command can generate a summary or detailed report of the analysis.

2. shell: This command creates a bash file containing all the commands executed via pyrpipe

3. benchmark: This command can generate benchmarks to compare walltimes of the different commands in the pipeline.

4. multiqc: This command uses the MultiQC tool to generate a report using pyrpipe logs and other logs geenrated by the pipeline.

## 2.3 Cookbook

**Contents**

- *Cookbook*
  - *Using SRA objects*
  - *Using RNASeqQC objects*
  - *Using Aligner objects*
  - *Using Assembler objects*
  - *Using Quantification objects*
  - *Using RNASeqTools objects*
  - *Using Runnable objects*
  - *Using pyrpipe_engine module*
  - *Using pyrpipe_utils module*

### 2.3.1 Using SRA objects

```python
from pyrpipe.sra import SRA #imports the SRA class

#create an SRA object using a valid run accession
"""
this checks if fastq files already exist in the directory,
otherwise downloads the fastq files and stores the path in the object
"""
myob=SRA('SRR1168424',directory='./results')

#create an SRA object using fastq paths
myob=SRA('SRR1168424',fastq='./fastq/SRR1168424_1.fastq',fastq2='./fastq/SRR1168424_2.
→fastq')

#create an SRA object using sra path
myob=SRA('SRR1168424',sra='./sra/SRR1168424.sra')

#accessing fastq files
print(myob.fastq,myob.fastq2)

#check if fastq files are present
print (myob.fastq_exists())

#check sra file
```

(continues on next page)

```
23  print (myob.sra_exists())
24
25  #delete fastq
26  myob.delete_fastq()
27
28  #delete sra
29  myob.delete_sra()
30
31  #download fastq
32  myob.download_fastq()
33
34  #trim fastq files
35  myob.trim(qcobject)
```

### 2.3.2 Using RNASeqQC objects

*RNASeqQC* objects can be used for quality control and trimming. These are defined in the *qc* module. Following example uses Trimgalore class but is applicable to any class extending the *RNASeqQC* class

```
1   from pyrpipe.qc import Trimgalore
2
3   #create trimgalore object
4   tgalore=Trimgalore()
5   #print category
6   print(tgalore._category) #should print Aligner
7
8   #use with SRA object
9   """
10  Following will trim fastq and update fastq paths in the sraobject
11  """
12  sraobject.trim(tgalore)
13  #following trims and returns qc
14  fq1,fq2=tgalore.perform_qc(sraobject)
15
16  #run trimgalore using user arguments; provide any arguments that Runnable.run() can␣
    ↪take
17  tgalore.run(*args,**kwargs)
```

### 2.3.3 Using Aligner objects

*Aligner* objects from the *mapping* module can be used to perform alignment tasks.

```
1   from pyrpipe.mapping import Star
2
3   #create a star object
4   star=Star(index='path_to_star_index')
5
6   #print category
7   print(star._category) #should print Aligner
8
9   #perform alignment using SRA object
10  bam=star.perform_alignment(sraobject)
11  #or
```

```
12  sraobject.align(star)
13  bam=sraobject.bam_path
14
15  #execute STAR with any arguments and parameters
16  kwargs={'--outFilterType' : 'BySJout',
17          '--runThreadN': '6',
18          '--outSAMtype': 'BAM SortedByCoordinate',
19          '--readFilesIn': 'SRR3098744_1.fastq SRR3098744_2.fastq'
20          }
21  star.run(**kwargs)
```

### 2.3.4 Using Assembler objects

*Assembler* objects are defined the the assembly module and can be used for transcript assembly.

```
1   from pyrpipe.assembly import Stringtie
2
3   #create a stringtie object
4   stringtie=Stringtie(guide='path_to_ref_gtf')
5
6   #perform assembly using SRA object
7   """
8   Note: following first runs star to perform alignment. After alignment the Sorted
9   BAM file is stored in the sraobject.bam_path attribute and returns the modified␣
    ↪sraobject
10  The assemble function requires a valid bam_path attribute to work.
11  """
12  sraobject.align(star).assemble(stringtie)
13  #Or manually set bam_path
14  sraobject.bam_path='/path/to/sorted.bam'
15  sraobject.assemble(stringtie)
16
17  #use perform_assembly function
18  result_gtf=stringtie.perform_assembly('/path/to/sorted.bam')
19
20  #run stringtie with user arguments
21  stringtie.run(verbose=True, **kwargs)
```

### 2.3.5 Using Quantification objects

*Quantification* type objects can perform quantification and are defined inside the *quant* module.

```
1   from pyrpipe.assembly import Salmon
2
3   #create salmon object
4   """
5   A valid salmon idex is required. If index is not found it is built using the provided␣
    ↪transcriptome
6   """
7   salmon=Salmon(index='path/to/index',transcriptome='path/to/tr')
8
9   #directly quantify using SRA object
10  sraobject.quant(salmon)
```

```python
11  #or trim reads before quant
12  sraobject.trim(tgalore).quant(salmon)
13  print('Result file',sraobject.abundance)
14
15  #use perform quant function
16  abundance_file=salmon.perform_quant(sraobject)
17
18  #use salmon with user defined arguments
19  salmon.run(**kwargs)
```

### 2.3.6 Using RNASeqTools objects

The *RNASeqTools* type is defined in *tools* module. This contains various tools used routinely for RNA-Seq data processing/analysis

```python
1   from pyrpipe.tools import Samtools
2
3   #create samtools object
4   samtools=Samtools(threads=6)
5
6   #convert sam to sorted bam
7   bam=samtools.sam_sorted_bam('sam_file')
8
9   #merge bam files
10  mergedbam=samtools.merge_bam(bamfiles_list)
11
12
13  #run samtools with used defined arguments
14  """
15  NOTE: the Runnable.run() method accepts a subcommand argument that allows user to
    ↪procide a subcommand like samtools index or samtools merge
16  """
17  samtools.run(*args,subcommand='index',**kwargs)
```

### 2.3.7 Using Runnable objects

The *Runnable* class, defined inside the *runnable* module, is the main parent class in *pyrpipe* i.e. all other classes borrows its functionality. User can directly create *Runnable* objects to define their own tools/commands.

A full example to build APIs is here: *Building APIs*

```python
1   from pyrpipe.runnable import Runnable
2
3   #say you want to use the tool orfipy
4   orfipy=Runnable(command='orfipy')
5
6   #execute orfipy as
7   orfipy.run(*args,**kwargs)
8
9   #another example using Unix grep
10  grep=Runnable(command='grep')
11  grep.run('query1','file1.txt',verbose=True)
12  grep.run('query2','file2.txt',verbose=True)
```

```
13
14   #extend Runnable to build more complex APIs that fit with each other
15   """
16   One can create classes extending the Runnable class.
17   Full example is given in the tutorial
18   """
19   class Orfipy(Runnable):
20       def __init__(self,*args,threads=None,mem=None,**kwargs):
21           super().__init__(*args,command='orfipy',**kwargs)
22           self._deps=[self._command]
23           self._param_yaml='orfipy.yaml'
24
25       #create special API functions that can work with other objects
26       def find_orfs(self,sra_object):
27           #define logic here and gather command options and parameters
28
29           #call the self.run() function and check values
30
31           #return a useful value
```

### 2.3.8 Using pyrrpipe_engine module

The *pyrrpipe_engine* module contain functions that creates new processes and enable executing commands. User can directly import *pyrrpipe_engine* module and start using the functions. This is very useful for quickly executing commands without having to create a Runnable object. A table describing functions implements in *pyrrpipe_engine* is provided in the tutorial *pyrrpipe_engine module*

```
1    import pyrrpipe_engine as pe
2
3    #execute_command: Runs a command, logs the status and returns the status (True or
     ↪False)
4    pe.execute_command(['ls', '-l'],logs=False,verbose=True)
5
6    #get_shell_output Runs a command and returns a tuple (returncode, stdout and stderr)
7    """
8    NOTE: only this function supports shell=True
9    """
10   result=pe.get_shell_output(['head','sample_file'])
11   #result contains return code, stdout, stderr
12   print(result)
13
14   #make a function dry-run compatible
15   """
16   when --dry-run flag is used this function will be skipped and the first parameter 'cmd
     ↪' will be printed to screen
17   """
18   @pe.dryable
19   def func(cmd,...)
20       #function logic here
21
22   """
23   Other way to work with dry-run flag is to directly import _dryrun flag
24   """
25
26   from from pyrrpipe import _threads,_force,_dryrun
```

```python
27  def myfunction(...):
28      if _dryrun:
29          print('This is a dry run')
30          return
31
32      #real code here...
```

### 2.3.9 Using pyrpipe_utils module

The *pyrpipe_utils* module contains several helpful functions that one needs to frequently access for typical computational pipelines. A table describing functions implements in *pyrpipe_utils* is provided in the tutorial *pyrpipe_utils module*

```python
1   import pyrpipe_utils as pu
2
3   #check if files exist
4   pu.check_files_exist('path/f1','path/f2','path/f3') #returns bool
5
6   #get filename without extention
7   pu.get_file_basename('path/to/file.ext')
8
9   #get file directory
10  pu.get_file_directory('path/to/file.ext')
11
12  #create a directory
13  pu.mkdir('path/to/dir')
14
15  #Search .txt files in a directly
16  pu.find_files('path/to/directory','.*\.txt$',recursive=False,verbose=False)
```

Contents:

- genindex

- modindex

- search