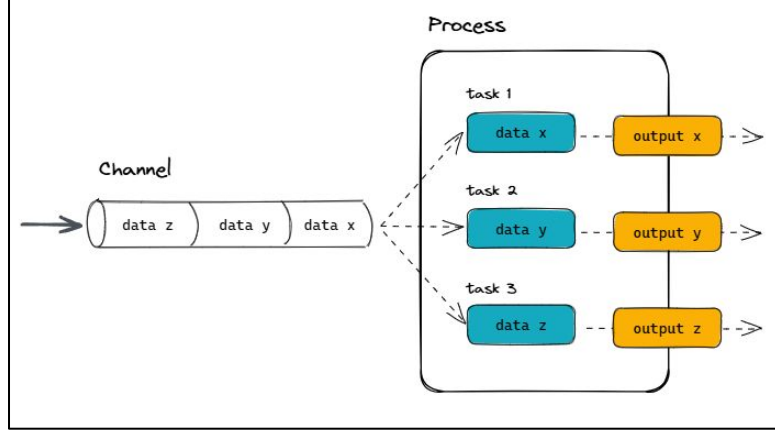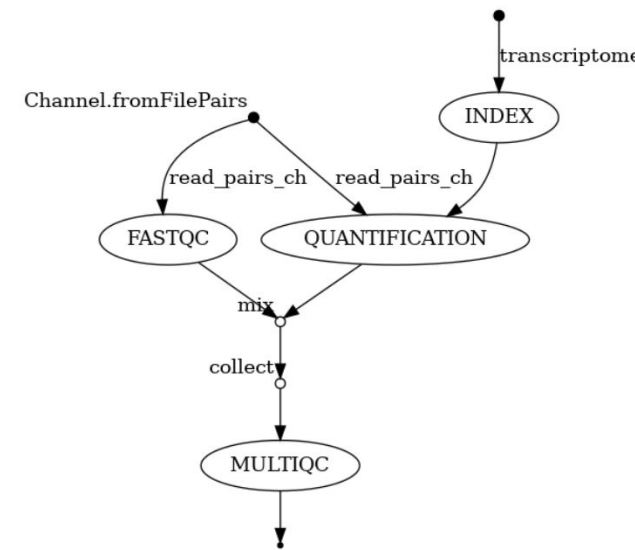## Simple RNA-Seq workflow

**Overview**

1. Indexes a transcriptome file
2. Performs quality controls
3. Performs quantification
4. Creates a MultiQC report

**Requirements**

Softwares: Nextflow-training dir in Gitpod
Languages: Python|Bash
Scripts: script7.nf
Tools: Salmon, FastQC, MultiQC

**Links**

Video tutorial
https://www.youtube.com/watch?v=nzR8DRq13oE

Written tutorial
https://training.nextflow.io/basic_training/rnaseq_pipeline/

Software link
https://nextflowio-training-2hwzxy3c0vu.ws-us117.gitpod.io/

---

## Definitions

A Nextflow workflow is made by joining processes.
Processes can be written in any Linux executable scripting language and are executed independently. Processes only communicate via asynchronous queues called channels.
Every input and output of a processes is a channel. The interaction between these processes are defined by input and output declarations.
Process is defined by input, output and script.

Nextflow pipeline represented as a direct acrylic graph (DAG). Vertices represent pipeline's processes and operators. Edges represent data dependencies (channels) between them.



nf-training -> nextflow run script7.nf -resume --reads 'data/ggal/"_{1,2}.fq' -with-dag flowchart.png



---

## Scripts

**script1.nf**

a. Define workflow parameters
   i. Inputs
      1. $projectDir
      2. Fastq's
      3. transcriptome.fa
   ii. Outputs
      1. **Multiqc outdir**
      2. **Results outdir**
   iii. Notes
      1. log.info command outputs multiline message and saves copy of info into log execution file
      2. $projectDir is a dynamic variable whose value is determined at runtime when and where the workflow is executed
      3. $projectDir path = "/workspace/gitpod/nf-training"

**script2.nf**

a. Create transcriptome index file
   i. Inputs
      1. transcriptome.fa path = params.transcriptome_file
   ii. Outputs
      1. Workflow scope containing index_ch channel
      2. Salmon_index outputs to index_ch #indexed transcriptome
   iii. Notes
      1. process INDEX() uses salmon
      2. Ensure nextflow.config defines salmon!!! See docker container image (line 1)
      3. Nextflow.config must be in $projectDir (currentdir)
      4. Ensure docker is enabled (line 3)

**script3.nf**

a. Collect read files by pairs into a channel
   i. Inputs
      1. params.reads #location of fq's
      2. fromFilePairs #known channel factory
   ii. Outputs
      1. read_pair_ch #a new channel variable
   iii. Notes
      1. fromFilePairs input is glob pattern and returns a channel of tuples. Each tuple contains two items: 1) read pair prefix 2) list of file paths
      2. .set ~ = , just used to assign
      3. Bonus points for wildcard function, but must be wrapped in single quotes if used

**script4.nf**

a. Expression quantification → with process QUANTIFICATION(read_pair_ch)
   i. Inputs
      1. Index_ch #previous output from process INDEX()
      2. read_pair_ch #previous output from known fromFilePairs channel
   ii. Outputs
      1. Workflow scope containing quant_ch channel
   iii. Notes
      1. Add tag to create more readable execution log
      2. Add publishDir to store the process results in outdir of choice (already defined as results/

**script5.nf**

a. FASTQC
   i. Inputs
      1. read_pair_ch #previous output from known fromFilePairs channel
   ii. Outputs
      1. Fqc's into sample_id folders
      2. fastqc_ch channel #contains fqc paths
   iii. Notes
      1. process FASTQC() makes a fastqc folder for each sample

**script6.nf**

a. MultiQC report
   i. Inputs
      1. fastqc_ch channel #contains fqc paths
   ii. Outputs
      1. final report in results/ in current workdir
   iii. Notes
      1. process MULTIQC()
      2. mix and collect operators together gather the outputs of quant_ch and fastq_ch as single input to ensure return of complete channel connects as single element

**script7.nf**

a. Handle event completion
   i. Notes
      1. Execute action following workflow completion (workflow.onComplete)
      2. Asynchronous tasks
      3. Bonus points for adding email notification in nextflow.config

---

## Define workflow parameters

```
OPEN EDITORS
  GROUP 1
    script7.nf nf-t...  2
  × script1.nf nf-t...  1
  GROUP 2
    nextflow.config nf...
  × script2.nf nf-tr...  1
GITPOD-WS (WORKSPACE)
  nf-training
   > .nextflow
   ∨ data
     ∨ ggal
       gut_1.fq
       gut_2.fq
       liver_1.fq
       liver_2.fq
       lung_1.fq
       lung_2.fq
       transcriptome.fa
   > index
```

```
nf-training > script1.nf
 1  #!/usr/bin/env nextflow
 2
 3  /*
 4  *pipeline input parameters
 5  */
 6
 7
 8  params.reads = "$projectDir/data/ggal/gut_{1,2}.fq"
 9  params.transcriptome_file = "$projectDir/data/ggal/transcriptome.fa"
10  params.multiqc = "$projectDir/multiqc"
11  params.outdir = "results"
12
13  log.info """\
14      R N A S E Q - N F   P I P E L I N E
15      ===================================
16      transcriptome: ${params.transcriptome_file}
17      reads        : ${params.reads}
18      outdir       : ${params.outdir}
19      """
20      .stripIndent()
```

## Process INDEX() & nextflow.config

```
24  process INDEX {
25      cpus 2
26      input:
27      path transcriptome
28
29      output:
30      path 'salmon_index'
31
32      script:
33      """
34      salmon index --threads $task.cpus -t $transcriptome -i salmon_index
35      """
36  }
37
   Preview DAG
38  workflow {
39      index_ch = INDEX(params.transcriptome_file)
40      index_ch.view()
41  }
```

```
nextflow.config  ×
nf-training > nextflow.config
 1  process.container = 'nextflow/rnaseq-nf'
 2  docker.runOptions = '-u $(id -u):$(id -g)'
 3  docker.enabled = true
```

## Collect paired fq files

```
   Preview DAG
84  workflow {
85      Channel
86          .fromFilePairs(params.reads, checkIfExists: true)
87          .set { read_pairs_ch }
88
89      index_ch = INDEX(params.transcriptome_file)
90
91
92
93  }
94
```

## Expression Quantification

```
36  process QUANTIFICATION {
37      tag "Salmon on $sample_id"
38      publishDir params.outdir, mode:'copy'
39
40      input:
41      path salmon_index                    → from index_ch #path to indexed transcriptome
42      tuple val(sample_id), path(reads)     → from read_pairs_ch #contains a tuple (value: sample_id/pattern, paths-to-reads: fqs)
43
44      output:
45      path "$sample_id"
46
47      script:
48      """
49      salmon quant --threads $task.cpus --libType=U -i $salmon_index -1 ${reads[0]} -2 ${reads[1]} -o $sample_id
50      """
   Preview DAG
84  workflow {
85      Channel
86          .fromFilePairs(params.reads, checkIfExists: true)
87          .set { read_pairs_ch }
88
89      index_ch = INDEX(params.transcriptome_file)
90      quant_ch = QUANTIFICATION(index_ch, read_pairs_ch)
91
92
93  }
94
```

## FASTQC, MULTIQC, Event completion

```
53  process FASTQC {
54      tag "FASTQC on $sample_id"
55
56      input:
57      tuple val(sample_id), path(reads)  →from read_pairs_ch #contains a tuple (value: sample_id/pattern, paths-to-reads: fqs)
58
59      output:
60      path "fastqc_${sample_id}_logs"
61
62      script:
63      """
64      mkdir fastqc_${sample_id}_logs
65      fastqc -o fastqc_${sample_id}_logs -f fastq -q ${reads}
66      """
67  }
68
69  process MULTIQC {
70      publishDir params.outdir, mode:'copy'
71
72      input:
73      path '*'
74
75      output:
76      path 'multiqc_report.html'
77
78      script:
79      """
80      multiqc .
81      """
82  }
   Preview DAG
84  workflow {
85      Channel
86          .fromFilePairs(params.reads, checkIfExists: true)
87          .set { read_pairs_ch }
88
89      index_ch = INDEX(params.transcriptome_file)
90      quant_ch = QUANTIFICATION(index_ch, read_pairs_ch)
91      fastqc_ch = FASTQC(read_pairs_ch)
92      MULTIQC(quant_ch.mix(fastqc_ch).collect())
93  }
94
95  workflow.onComplete {
96      log.info ( workflow.success ? "\nDone! Open the following report in your browser --> $params.outdir/multiqc_report.html\n" : "Oops .. something went wrong" )
97  }
98
```

---

## Additional notes

1. **Project Structure:**
   - my_workflow/
     - main.nf
     - data/
     - scripts/
       - script1.nf
       - script2.nf

2. **If you are in the scripts/ directory:**
   - $projectDir will refer to the my_workflow/ directory, not the scripts/ directory.