

Using configuration YAML files as a one-stop for all custom parameter needs

BYOB 06-08-2022

Caroline Esnault

NICHD Bioinformatics and Scientific Programming Core

Snakefile-hc ?

```
1 import sys
2 import os
3
4
5
6
7
8
9
10
11 # define targets to create
12 final_targets = ['data/featurecounts.tsv']
13
14 # uncomment line below to run differential expression
15 final_targets.append('downstream/rnaseq.html')
16
17 rule targets:
18     """
19     Final targets to create
20     """
21     input: final_targets
22
23
24 # featurecounts
25 rule featurecounts:
26     """
27     Count reads in annotations with featureCounts from the subread package
28     """
29     input:
30         annotation='../raw-data/ecoli.gtf',
31         bam='../raw-data/example.bam'
32     params:
33         strand='s2'      # change here to match library strandedness
34     output:
35         counts='data/featurecounts.txt'
36     log:
37         'logs/featurecounts.txt.log'
38     shell:
39         'featureCounts '
40         '-{params.strand} -t exon -g derived_id ' # change to specify columns
41         '-p '
42         '-a {input.annotation} '
43         '-o {output.counts} '
44         '{input.bam} '
45         '&> {log}'
46
```

What should *not* be done

comment out / uncomment
depending on needs

hard-coded paths

hardcoded parameters to modify

```

1 #!/usr/bin/perl -w
2
3 ## Script to truncate the last 50bp of the sequences from sequence file
4 ## (MiSeq for example).
5 ## The goal is to have the same structure of sequences between HiSeq
6 ## (101bp run) and MiSeq (151bp run)
7
8
9 my $n = 1;
10 my $seqline;
11 my $lenght = 73;      # can modify here the lenght of sequence to keep
12
13 # ask for the sequence file
14 print "fastq file:";
15 chomp($fastqfile = <STDIN>);
16
17 $output = $fastqfile;
18 $output =~ s/.txt/_truncated73.txt/;
19 $output =~ s/.fastq/_truncated73.txt/;
20
21 open OUTPUT, ">$output" or die "can't open '$output': $!";
22
23 open (FASTQFILE, $fastqfile);
24 print "Screening the sequences\n";
25 while (<FASTQFILE>) {
26     $seqline = $_;
27     chomp $seqline;
28     # truncate at the first 101bp if it is a sequence lines (1 / 4)
29     if ($n == 2) {
30         $seqline = substr ($seqline, 0, $lenght);...
31     }
32     if ($n == 4) {
33         $n = 0;
34     }
35     print OUTPUT "$seqline\n";
36     $n += 1;
37 }
38
39 close(FASTQFILE);
40
41
42 close OUTPUT;
43
44
45
46

```

Real-life example of what not to do

Goals

- move all parameters and variables OUTSIDE of the code
- no hardcoding in script
- main version of code fits all experiments
- easy 1-stop way to find which parameters were used

YAML language

definition:

YAML (rhymes with “camel”) is a data serialization language designed to be human-friendly and work well with modern programming languages for common everyday tasks.

Y : Yet	Y : YAML
A : Another	A : Ain't
M : Markup	M : Markup
L : Language	L : Language
(source: redhat.com)	recursive acronym (source yaml.org)

timeline:

- **2001**: 1st YAML framework written in Perl
- **2003**: Ruby 1st language to ship a YAML framework as part of its core language
- **2004**: YAML 1.0 specification published by Clark Evans, Oren Ben-Kiki, and Ingy döt Net
- **2005**: YAML 1.1
- **2006**: Kyrlo Simonov produced PyYAML and LibYAML
- **2009**: YAML 1.2
- **2020**: new YAML language design team began meeting regularly to discuss improvements
- **2021**: YAML 1.2.2

```

config.yaml ?
1 # YAML
2 ---
3
4 # comment
5 key: some_text
6
7 some_sequence: ['item1', 'item2']
8
9 another_sequence:
10   _ 'item1'
11   _ 'item2'
12
13 some_mapping: {'control': 'blue', 'mutant': 'green'}
14
15 another_mapping:
16   'control': 'blue'
17   'mutant': 'green'
18
19 mapping_of_mappings:
20   mapping_of_sequence:
21     _ level1
22     _ level2
23     _ level3
24   nested_mapping_of_mappings:
25     'red'
26     'orange'
27     'yellow'
28
29 plain:
30   This unquoted scalar
31   spans many lines.
32
33 quoted:
34   "So does this
35   quoted scalar.\n"
36
37

```

YAML structure

for comments

three dashes (“---”) to separate directives from document content

scalar

sequence (ordered)

mapping (not ordered)

single or double quote, or even no quote

structure is determined by *indentation* (not TAB!)

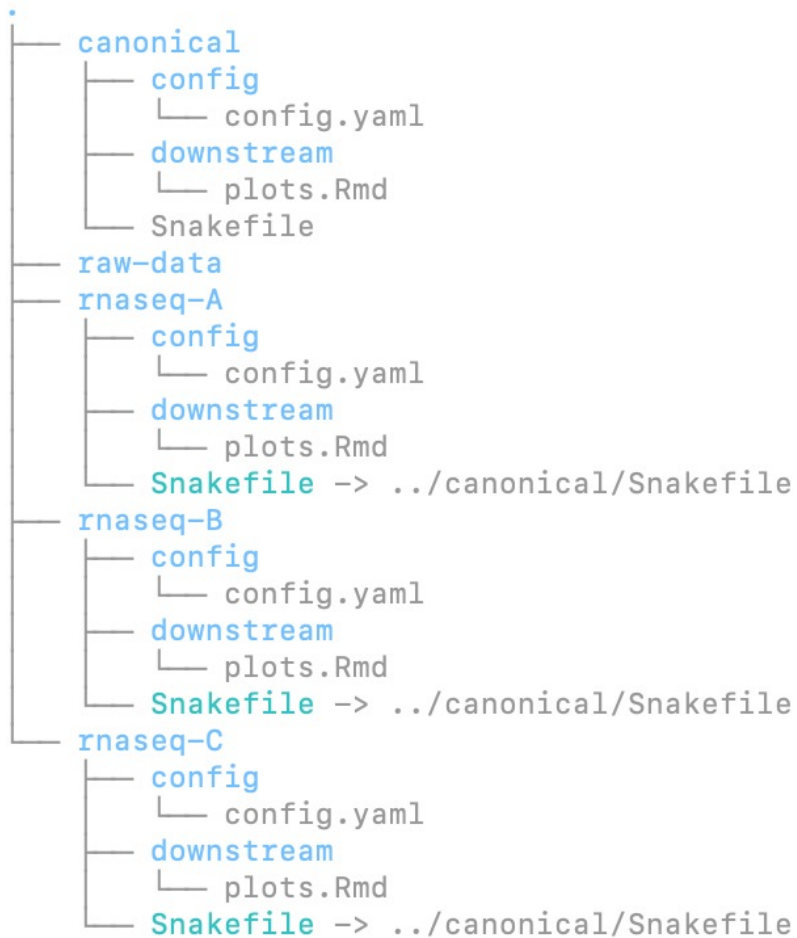
scalars (quoted or unquoted) can span many lines

Parsers exist in many languages

- perl

```
use strict;  
use warnings;  
use YAML::XS 'LoadFile';  
use Data::Dumper;  
  
my $config = LoadFile('config.yaml');  
  
print Dumper($config);
```

- Java
- Javascript
- C/C++
- Rust
- Shell
- Python
- R
- ...
- anything? structure makes it possible to parse relatively easily



Use unique script in complex project structure

main copy in canonical

change 1 copy to change all

can use the same config.yaml for python and downstream Rmd

Loading YAML in python

```
1 import sys
2 import os
3 import yaml
4
5
6 # load config
7 configfn= 'config/config.yaml'
8 config = yaml.safe_load(open(configfn))
9
10
```

Snakefile-hc ?

```

1 import sys
2 import os
3
4
5
6
7
8
9
10
11 # define targets to create
12 final_targets = ['data/featurecounts.tsv']
13
14 # uncomment line below to run differential expression
15 final_targets.append('downstream/rnaseq.html')
16
17 rule targets:
18     """
19     Final targets to create
20     """
21     input: final_targets
22
23
24 # featurecounts
25 rule featurecounts:
26     """
27     Count reads in annotations with featureCounts from the subread package
28     """
29     input:
30         annotation='../raw-data/ecoli.gtf',
31         bam='../raw-data/example.bam'
32     params:
33         strand='s2'      # change here to match library strandedness
34     output:
35         counts='data/featurecounts.txt'
36     log:
37         'logs/featurecounts.txt.log'
38     shell:
39         'featureCounts '
40         '-{params.strand} -t exon -g derived_id ' # change to specify columns
41         '-p '
42         '-a {input.annotation} '
43         '-o {output.counts} '
44         '{input.bam} '
45         '&> {log}'
46

```

config.yaml ?

```

1 # rnaseq-A YAML
2 ---
3
4 # Snakefile parameters
5 # -----
6
7 # annotation
8 gtf: ../raw-data/ecoli.gtf
9
10 # BAM file
11 bam: ../raw-data/example.bam
12
13 # Extra featureCounts parameters
14 extra: '-s2 -t exon -g derived_id'
15
16 # output directory
17 outdir: data
18
19 # True or False to run differential expression
20 differential_expression: True
21
22
23 # downstream DE Rmd parameters
24 # -----
25
26 # title
27 title: 'Differential expression RNAseq A'
28
29 # sampletable
30 sampletable: '../config/sampletable.tsv'
31
32 # featurecounts file
33 featurecounts: '../data/featurecounts.txt'
34
35 # reference level for dds object
36 level.column: 'group'
37 level.ref: 'log'
38
39 # contrast(s): indicate for each contrast the ctrst.group,
40 # ctrst.test, ctrst.control
41 contrasts:
42     mutant_vs_ctrl:
43         ctrst.group: 'group'
44         ctrst.test: 'exp'
45         ctrst.control: 'log'
46
47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49     - 'G0-9381'
50     - 'EG10747'
51

```

```

Snakefile [?] bash [?]
1 import sys
2 import os
3 import yaml
4
5
6 # load config
7 configfn= 'config/config.yaml'
8 config = yaml.safe_load(open(configfn))
9
10
11 # define targets to create
12 final_targets = [config['outdir'] + '/featurecounts.tsv']
13
14 if config['differential_expression']:
15     final_targets.append('downstream/rnaseq.html')
16
17 rule targets:
18     """
19     Final targets to create
20     """
21     input: final_targets
22
23
24 # featurecounts
25 rule featurecounts:
26     """
27     Count reads in annotations with featureCounts from the subread package
28     """
29     input:
30         annotation=config['gtf'],
31         bam=config['bam']
32     params:
33         extra=config['extra']
34     output:
35         counts=config['outdir'] + '/featurecounts.txt'
36     log:
37         'logs/featurecounts.txt.log'
38     shell:
39         'featureCounts '
40         '{params.extra} '
41         '-p '
42         '-a {input.annotation} '
43         '-o {output.counts} '
44         '{input.bam} '
45         '&> {log}'
46

```

```

config.yaml [?]
1 # rnaseq-A YAML
2 ---
3
4 # Snakefile parameters
5 # -----
6
7 # annotation
8 gtf: ../raw-data/ecoli.gtf
9
10 # BAM file
11 bam: ../raw-data/example.bam
12
13 # Extra featureCounts parameters
14 extra: '-s2 -t exon -g derived_id'
15
16 # output directory
17 outdir: data
18
19 # True or False to run differential expression
20 differential_expression: True
21
22
23 # downstream DE Rmd parameters
24 # -----
25
26 # title
27 title: 'Differential expression RNAseq A'
28
29 # sampletable
30 sampletable: '../config/sampletable.tsv'
31
32 # featurecounts file
33 featurecounts: '../data/featurecounts.txt'
34
35 # reference level for dds object
36 level.column: 'group'
37 level.ref: 'log'
38
39 # contrast(s): indicate for each contrast the ctrst.group,
40 # ctrst.test, ctrst.control
41 contrasts:
42     mutant_vs_ctrl:
43         ctrst.group: 'group'
44         ctrst.test: 'exp'
45         ctrst.control: 'log'
46
47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49     - 'G0-9381'
50     - 'EG10747'
51

```

```

Snakefile [?] bash [?]
1 import sys
2 import os
3 import yaml
4
5
6 # load config
7 configfn= 'config/config.yaml'
8 config = yaml.safe_load(open(configfn))
9
10
11 # define targets to create
12 final_targets = [config['outdir'] + '/featurecounts.tsv']
13
14 if config['differential_expression']:
15     final_targets.append('downstream/rnaseq.html')
16
17 rule targets:
18     """
19     Final targets to create
20     """
21     input: final_targets
22
23
24 # featurecounts
25 rule featurecounts:
26     """
27     Count reads in annotations with featureCounts from the subread package
28     """
29     input:
30         annotation=config['gtf'],
31         bam=config['bam']
32     params:
33         extra=config['extra']
34     output:
35         counts=config['outdir'] + '/featurecounts.txt'
36     log:
37         'logs/featurecounts.txt.log'
38     shell:
39         'featureCounts '
40         '{params.extra} '
41         '-p '
42         '-a {input.annotation} '
43         '-o {output.counts} '
44         '{input.bam} '
45         '&> {log}'
46

```

```

config.yaml [?]
1 # rnaseq-A YAML
2 ---
3
4 # Snakefile parameters
5 # -----
6
7 # annotation
8 gtf: ../raw-data/ecoli.gtf
9
10 # BAM file
11 bam: ../raw-data/example.bam
12
13 # Extra featureCounts parameters
14 extra: '-s2 -t exon -g derived_id'
15
16 # output directory
17 outdir: data
18
19 # True or False to run differential expression
20 differential_expression: True
21
22
23 # downstream DE Rmd parameters
24 # -----
25
26 # title
27 title: 'Differential expression RNAseq A'
28
29 # sampletable
30 sample:
31
32 # featureCounts
33 featureCounts:
34
35 # reference
36 level: 1
37 level: 1
38
39 # contrasts
40 # contrast
41 contrasts:
42     mutant_vs_ctrl:
43         ctrst.group: 'group'
44         ctrst.test: 'exp'
45         ctrst.control: 'log'
46
47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49     - 'G0-9381'
50     - 'EG10747'
51

```

```

In [21]: config
Out[21]:
{'gtf': '../raw-data/ecoli.gtf',
 'bam': '../raw-data/example.bam',
 'extra': '-s2 -t exon -g derived_id',
 'outdir': 'data',
 'differential_expression': True}

```

Loading YAML in R

```
34  
35 library(yaml)  
36 cfgfn <- '../config/congif.yaml'  
37 cfg <- read_yaml(cfgfn)  
38  
39
```

```

plots.Rmd [?] bash [?]
1 ---
2 output:
3   html_document:
4     code_folding: hide
5     toc: true
6     toc_float: true
7     toc_depth: 3
8 ---
9
10
11 ```{r}
12 knitr::opts_chunk$set(message=FALSE, warning=FALSE)
13 ```
14
15 ```{r libraries}
16 library(ggplot2)
17 library(tidyr)
18 library(dplyr)
19 library(DESeq2)
20 library(yaml)
21 ```
22
23 ```{r configs}
24 # read config yaml
25 cfg <- read_yaml('../config/config.yaml')
26 ```
27
28 # `r cfg[['title']]`
29
30 ```{r coldata_setup, cache=TRUE}
31 # Set up all of the metadata for the samples and experimental design.
32 colData <- read.table(cfg[['samletable']], sep='\t', header=TRUE, stringsAsFactors=FALSE)
33
34 # releve factor
35 level.col <- cfg[['level.column']]
36 level.ref <- cfg[['level.ref']]
37 colData[[level.col]] <- as.factor(colData[[level.col]])
38 colData[[level.col]] <- releve(colData[[level.col]], ref=level.ref)
39 rownames(colData) <- colData[,1]
40
41 knitr::kable(colData)
42 ```
43
44 # DESeq Data Set
45
46 ```{r dds}
47 # count matrix from featureCounts
48 cts <- read.table(cfg[['featurecounts']], sep='\t', header=TRUE, row.names='Geneid')
49 cts <- cts[colData[['samlename']]]
50
51 # dds
52 dds <- DESeqDataSetFromMatrix(countData = cts,
53                               colData = colData,
54                               design = formula(paste("~", level.col)))
55 dds <- DESeq(dds)
56 print(dds)
57 ```
58
59 # Contrasts results

```

```

config.yaml [?]
1 # rnaseq-A YAML
2 ---
3
4 # Snakefile parameters
5 # -----
6
7 # annotation
8 gtf: ../raw-data/ecoli.gtf
9
10 # BAM file
11 bam: ../raw-data/example.bam
12
13 # Extra featureCounts parameters
14 extra: '-s2 -t exon -g derived_id'
15
16 # output directory
17 outdir: data
18
19 # True of False to run differential expression
20 differential_expression: True
21
22
23 # downstream DE Rmd parameters
24 # -----
25
26 # title
27 title: 'Differential expression RNAseq A'
28
29 # samletable
30 samletable: '../config/samletable.tsv'
31
32 # featurecounts file
33 featurecounts: '../data/featurecounts.txt'
34
35 # reference level for dds object
36 level.column: 'group'
37 level.ref: 'log'
38
39 # contrast(s): indicate for each contrast the ctrst.group,
40 # ctrst.test, ctrst.control
41 contrasts:
42   mutant_vs_ctrl:
43     ctrst.group: 'group'
44     ctrst.test: 'exp'
45     ctrst.control: 'log'
46
47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49   _ 'G0-9381'
50   _ 'EG10747'
51

```



```

plots.Rmd [?] bash [?]
1 ---
2 output:
3   html_document:
4     code_folding: hide
5     toc: true
6     toc_float: true
7     toc_depth: 3
8 ---
9
10
11 ```{r}
12 knitr::opts_chunk$set(message=FALSE, warning=FALSE)
13 ```
14
15 ```{r libraries}
16 library(ggplot2)
17 library(tidyr)
18 library(dplyr)
19 library(DESeq2)
20 library(yaml)
21 ```
22
23 ```{r configs}
24 # read config yaml
25 cfg <- read_yaml('../config/config.yaml')
26 ```
27
28 # `r cfg[['title']]`
29
30 ```{r coldata_setup, cache=TRUE}
31 # Set up all of the metadata for the samples and experimental design.
32 colData <- read.table(cfg[['samletable']], sep='\t', header=TRUE, stringsAsFactors=FALSE)
33
34 # releve factor
35 level.col <- cfg[['level.column']]
36 level.ref <- cfg[['level.ref']]
37 colData[[level.col]] <- as.factor(colData[[level.col]])
38 colData[[level.col]] <- relevel(colData[[level.col]], ref=level.ref)
39 rownames(colData) <- colData[,1]
40
41 knitr::kable(colData)
42 ```
43
44 # DESeq Data Set
45
46 ```{r dds}
47 # count matrix from featureCounts
48 cts <- read.table(cfg[['featurecounts']], sep='\t', header=TRUE, row.names='Geneid')
49 cts <- cts[colData[['samplename']]]
50
51 # dds
52 dds <- DESeqDataSetFromMatrix(countData = cts,
53                               colData = colData,
54                               design = formula(paste("~", level.col)))
55 dds <- DESeq(dds)
56 print(dds)
57 ```
58
59 # Contrasts results

```

```

config.yaml [?]
1 # rnaseq-A YAML
2 ---
3
4 # Snakefile parameters
5 # -----
6
7 # annotation
8 gtf: ../raw-data/ecoli.gtf
9
10 # BAM file
11 bam: ../raw-data/example.bam
12
13 # Extra featureCounts parameters
14 extra: '-s2 -t exon -g derived_id'
15
16 # output directory
17 outdir: data
18
19 # True or False to run differential expression
20 differential_expression: True
21
22
23 # downstream DE Rmd parameters
24 # -----
25
26 # title
27 title: 'Differential expression RNAseq A'
28
29 # samletable
30 samletable: '../config/samletable.tsv'
31
32 # featurecounts file
33 featurecounts: '../data/featurecounts.txt'
34
35 # reference level for dds object
36 level.column: 'group'
37 level.ref: 'log'
38
39 # contrast(s): indicate for each contrast the d
40 # cntrst.test, cntrst.control
41 contrasts:
42   mutant_vs_ctrl:
43     cntrst.group: 'group'
44     cntrst.test: 'exp'
45     cntrst.control: 'log'
46
47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49   _ 'G0-9381'
50   _ 'EG10747'
51

```

```

> cfg
$gtf
[1] "../raw-data/ecoli.gtf"

$bam
[1] "../raw-data/example.bam"

$extra
[1] "-s2 -t exon -g derived_id"

$outdir
[1] "data"

$differential_expression
[1] TRUE

$title
[1] "Differential expression RNAseq A"

$samletable
[1] "../config/samletable.tsv"

$featurecounts
[1] "../data/featurecounts.txt"

$level.column
[1] "group"

$level.ref
[1] "log"

$contrasts
$contrasts$mutant_vs_ctrl
$contrasts$mutant_vs_ctrl$cntrst.group
[1] "group"

$contrasts$mutant_vs_ctrl$cntrst.test
[1] "exp"

$contrasts$mutant_vs_ctrl$cntrst.control
[1] "log"

```

Loop through contrasts

plots.Rmd

```
59 # Contrasts results
60
61 ```{r contrasts}
62 # results
63 res <- list()
64 for (ctrst in names(cfg[['contrasts']])) {
65   res[[ctrst]] <- lfcShrink(dds,
66                             contrast=c(cfg[['contrasts']][[ctrst]][['ctrst.group']],
67                                         cfg[['contrasts']][[ctrst]][['ctrst.test']],
68                                         cfg[['contrasts']][[ctrst]][['ctrst.control']],
69                                         type='normal')
70   cat(paste0('\n\n## ', ctrst, '\n\n'))
71   print(res[[ctrst]])
72 }
73 ```
```

config.yaml

```
38
39 # contrast(s): indicate for each contrast the ctrst.group,
40 # ctrst.test, ctrst.control
41 contrasts:
42   mutant_vs_ctrl:
43     ctrst.group: 'group'
44     ctrst.test: 'exp'
45     ctrst.control: 'log'
46
```


Define conditional rule / chunk

plots.Rmd

```

75 # Conditional plotting
76
77 ```{r}
78 if ('genes.to.plot' %in% names(cfg)) {
79   eval_plot <- TRUE
80 } else {
81   eval_plot <- FALSE
82   cat('\n\nNot run\n\n')
83 }
84 ```
85
86 ```{r plots, eval=eval_plot}
87 for (gene in cfg[['genes.to.plot']]) {
88   plotCounts(dds, gene=gene, intgroup=cfg[['contrasts']][[ctrst]][['ctrst.group']])
89 }
90 ```
91

```

config.yaml

```

47 # genes of interest to plot, comment out if none
48 genes.to.plot:
49   - 'G0-9381'
50   - 'EG10747'
51

```

Snakefile

```

13
14 if config['differential_expression']:
15   final_targets.append('downstream/rnaseq.html')
16
48 # conditional rule: downstream DESeq2 analysis
49 if config['differential_expression']:
50
51   rule rnaseq_rmarkdown:
52     """
53     Run and render the RMarkdown file that performs differential expression
54     """
55     input:
56       featurecounts=rules.featurecounts.output,
57       rmd='downstream/rnaseq.Rmd',
58       sampletable=config['sampletable']
59     output:
60       'downstream/rnaseq.html'
61     log:
62       'downstream/rnaseq.log'
63     shell:
64       'Rscript -e '
65       '""rmarkdown::render("{input.rmd}")" '
66       '> {log} 2>&1'
67

```

config.yaml

```

18
19 # True or False to run differential expression
20 differential_expression: True
21
22

```

Get history of parameters when using version control

```
commit daf9178f66a691f54901ba940d268fd8e5660078
Author: Caroline Esnault <caroline.esnault@nih.gov>
Date: Wed Jun 8 21:35:03 2022 -0400

    rename contrast

diff --git a/rnaseq-A/config/config.yaml b/rnaseq-A/config/config.yaml
index 669d7b8..b636fbd 100644
--- a/rnaseq-A/config/config.yaml
+++ b/rnaseq-A/config/config.yaml
@@ -34,12 +34,12 @@ featurecounts: '../data/featurecounts.txt'

# reference level for dds object
level.column: 'group'
-level.ref: 'exp'
+level.ref: 'log'

# contrast(s): indicate for each contrast the ctrst.group,
# ctrst.test, ctrst.control
contrasts:
- mutant_vs_ctrl:
+ exp_vs_log:
    ctrst.group: 'group'
    ctrst.test: 'exp'
    ctrst.control: 'log'

commit e19d6c55e3ea915f7cc2f1070b1156adf3898baf
Author: Caroline Esnault <caroline.esnault@nih.gov>
Date: Wed Jun 8 21:33:37 2022 -0400

    change level

diff --git a/rnaseq-A/config/config.yaml b/rnaseq-A/config/config.yaml
index c1aa100..669d7b8 100644
--- a/rnaseq-A/config/config.yaml
+++ b/rnaseq-A/config/config.yaml
@@ -34,7 +34,7 @@ featurecounts: '../data/featurecounts.txt'

# reference level for dds object
level.column: 'group'
-level.ref: 'log'
+level.ref: 'exp'

# contrast(s): indicate for each contrast the ctrst.group,
# ctrst.test, ctrst.control
```

```
git log -p -- rnaseq-A/config/config.yaml
```



Thanks for listening!