

*THE*

# RNAseq Analysis via Hisat, StringTie and Ballgown

*MANUAL*

Timo Lassmann



TODAY

# Contents

# Introduction

Included are two RNAseq workflows implemented in `bpipe`<sup>sadedin2012</sup>. Each produces a Ballgown<sup>frazee2014</sup>`flexible` R object for differential gene expression analysis.

# Quick Usage

To run these pipelines simply point them to the input fastq files:

```
bpipe run -n 4 -r ../pipelines/StringTieDeNovo.groovy *R*.fq.gz
```

Note: you may have to modify the regular expressions at the start of the pipelines to recognize your file names.

## 2.1 Workflow 1: Finding novel isoforms:

This workflow:

1. aligns all reads to the genome using HiSat<sup>kim2015hisat</sup>.
2. assembles mapped reads into transcripts using StringTie<sup>pertea2015stringtie</sup>. Known transcripts from gencode v19 are provided to guide the assembly.
3. merges transcripts from all samples using cuffmerge.
4. annotates transcripts using cuffcompare
5. quantifies transcripts from step 3 across all samples using StringTie<sup>pertea2015stringtie</sup>.

**File: StringTieDeNovoHUMAN.groovy**

```
##Settings

HISATINDEX=HISATHG19

GENOME=GENOMEHG19

KNOWNTRANSCRIPTS=GencodeV19_HS
KNOWNSPICESITES=GencodeV19_HS_SPLICESITES
```

```

MERGED_TRANSCRIPTS = "merged_asm/merged.gtf"

FINAL_TRANSCRIPTS = "merged_asm/finalmodel.combined.gtf"

Bpipe.run {"%_R*.fastq.gz" * [ hisat_align + stringtie] +
    makeassemblylist + cuffmerge + cuffcompare + "%_R*.fastq.gz" *
    [ hisat_align + stringtieB] + make_ballgown_obj}

```

## File: StringTieDeNovoMOUSE.groovy

```

##Settings

HISATINDEX=HISATMM10
GENOME=GENOMEMM10

KNOWNTRANSCRIPTS=GENCODEVM5_MM
KNOWNSPICESITES=GENCODEVM5_MM_SPLICESITES

MERGED_TRANSCRIPTS = "merged_asm/merged.gtf"

FINAL_TRANSCRIPTS = "merged_asm/finalmodel.combined.gtf"

Bpipe.run {"%_R*.fastq.gz" * [ hisat_align + stringtie] +
    makeassemblylist + cuffmerge + cuffcompare + "%_R*.fastq.gz" *
    [ hisat_align + stringtieB] + make_ballgown_obj}

```

## 2.2 Workflow 2: Quantification of known transcripts.

This workflow:

1. aligns all reads to the genome using HiSat<sup>kim2015hisat</sup>.
2. quantifies gencode v19 transcripts from across all samples using StringTie<sup>pertea2015stringtie</sup>.

## 2.3 Settings

```

load "localprograms.groovy"
load "localdatafiles.groovy"
load "../modules/RNaseq.groovy"

```

## File: make\_\_ballgown\_obj.sh

This is a small script to turn stringtie tables into an R object.

```
#!/bin/bash
```

```

echo "library(ballgown)" > ballgownR_glue.R
echo -ne "sample_directories <- c(" >> ballgownR_glue.R
counter=1
for line in $(find . -name "e2t.ctab"); do
  a="$(echo $line | rev | cut -d"/" -f2,3 | rev)"
  if [ $counter = 1 ]; then
    echo -ne "'$a'" >> ballgownR_glue.R
  else
    echo -ne ", '$a'" >> ballgownR_glue.R
  fi
  counter=$((counter + 1));
done
echo ")" >> ballgownR_glue.R
echo "bg <- ballgown(samples = sample_directories)" >>
  ballgownR_glue.R

echo "save(bg, file = 'bg.rda')" >> ballgownR_glue.R

Rscript ballgownR_glue.R

```