

Differential Expression Analysis

Transcriptome Analysis Overview

- 1) Download reads from sequencing facility or SRA
- 2) Check the quality of raw sequence files (FastQC)
- 3) Remove adapter sequence (e.g. CutAdapt)
- 4) Align RNAseq reads to the genome (e.g. STAR, HISAT2)
- 5) Get transcript abundances (e.g. HTSeq2)
- 6) QC individual count files
- 7) Make an expression matrix combining multiple samples into 1 file
- 8) QC matrix
- 9) Conduct biological analysis (e.g. determine differentially expressed genes with DESeq2)

Maize Gene Atlas

the plant journal



The Plant Journal (2011)

doi: 10.1111/j.1365-3113X.2011.04527.x

Genome-wide atlas of transcription during maize development

Rajandeep S. Sekhon^{1,2,*}, Haining Lin^{3,4,1}, Kevin L. Childs^{3,4}, Candice N. Hansey^{3,4}, C. Robin Buell^{3,4}, Natalia de Leon^{1,2} and Shawn M. Kaeppler^{1,2,*}

¹Department of Energy Great Lakes Bioenergy Research Center, University of Wisconsin-Madison, 1575 Linden Drive, Madison, WI 53706, USA,

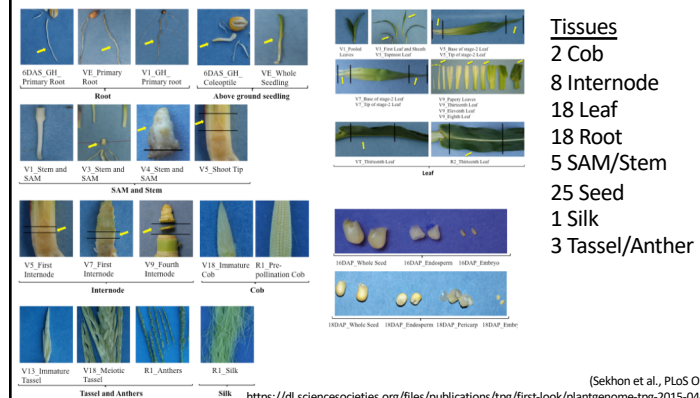
²Department of Agronomy, University of Wisconsin-Madison, 1575 Linden Drive, Madison, WI 53706, USA,

³Department of Energy Great Lakes Bioenergy Research Center, Michigan State University, East Lansing, MI 48824, USA, and

⁴Department of Plant Biology, Michigan State University, East Lansing, MI 48824, USA

Maize Gene Atlas

Catalog of gene expression throughout tissues throughout development



1) Download reads from SRA to MSI

```
#!/bin/bash
#SBATCH --time=3:00:00
#SBATCH --ntasks=8
#SBATCH --mem=16g
#SBATCH --tmp=16g
#SBATCH --mail-type=ALL
#SBATCH --mail-user=cnhirsch@um.edu

echo "Step 0: Setting up directory structure..."
mkdir /home/agro5431/cnhirsch/transcriptome_analysis/
cd /home/agro5431/cnhirsch/transcriptome_analysis/

echo "Step 1: Downloading the file from the SRA..."
module load sratoolkit/2.8.2
fastq-dump --split-files -f Q33 SRR948252

echo "Step 2: Check the quality of the raw reads..."
module load fastqc/0.11.7
fastqc -f fastq SRR948252.1.fastq

echo "Step 3: Remove the adapter sequences..."
module load cutadapt/1.18
cutadapt -a GATCGGAAGAGCACACGTCTGAACTCCAGTCACGATCATCTGATCGCTTCCTTCGTTG -a AATGATACGGCAGCAGACGATCTACATCTTCCCTACACGACGCTTCGTCATC -f fastq -m 30 --quality-base=33 -o SRR948252_1_cutadapt.fastq SRR948252.1.fastq > SRR948252_cutadapt.log

echo "Step 4: Check the quality of the cleaned reads..."
fastqc -f fastq SRR948252_1_cutadapt.fastq

echo "Step 5: Align RNAseq reads to the genome..."
module load star/2.5.3a
STAR --genomeDir /home/agro5431/shared/B73_STAR \
--runThreadN 8 \
--readFilesIn SRR948252_1_cutadapt.fastq \
--outFileNamePrefix SRR948252_1_cutadapt_STAR \
--outSAMtype BAM SortedByCoordinate \
--outSAMunmapped Within \
--outSAMattributes Standard

echo "Step 6: Get transcript abundance..."
module load htseq/0.7.2
htseq-count -f bam -s no -t gene -i ID -m union -a 20 SRR948252_1_cutadapt_STARAligned.sortedByCoord.out.bam /home/agro5431/shared/Zea_may6.AGPv4.33.gff3 > htseq_SRR948252.txt
```

2) Check the quality of raw sequence files

```
#!/bin/bash
#SBATCH --time=3:00:00
#SBATCH --ntasks=8
#SBATCH --mem=16g
#SBATCH --tmp=16g
#SBATCH --mail-type=ALL
#SBATCH --mail-user=cnhirsch@um.edu

echo "Step 0: Setting up directory structure..."
mkdir /home/agro5431/cnhirsch/transcriptome_analysis/
cd /home/agro5431/cnhirsch/transcriptome_analysis/

echo "Step 1: Downloading the file from the SRA..."
module load sratoolkit/2.8.2
fastq-dump --split-files -f Q33 SRR948252

echo "Step 2: Check the quality of the raw reads..."
module load fastqc/0.11.7
fastqc -f fastq SRR948252.1.fastq

echo "Step 3: Remove the adapter sequences..."
module load cutadapt/1.18
cutadapt -a GATCGGAAGAGCACACGTCTGAACTCCAGTCACGATCATCTGATCGCTTCCTTCGTTG -a AATGATACGGCAGCAGACGATCTACATCTTCCCTACACGACGCTTCGTCATC -f fastq -m 30 --quality-base=33 -o SRR948252_1_cutadapt.fastq SRR948252.1.fastq > SRR948252_cutadapt.log

echo "Step 4: Check the quality of the cleaned reads..."
fastqc -f fastq SRR948252_1_cutadapt.fastq

echo "Step 5: Align RNAseq reads to the genome..."
module load star/2.5.3a
STAR --genomeDir /home/agro5431/shared/B73_STAR \
--runThreadN 8 \
--readFilesIn SRR948252_1_cutadapt.fastq \
--outFileNamePrefix SRR948252_1_cutadapt_STAR \
--outSAMtype BAM SortedByCoordinate \
--outSAMunmapped Within \
--outSAMattributes Standard

echo "Step 6: Get transcript abundance..."
module load htseq/0.7.2
htseq-count -f bam -s no -t gene -i ID -m union -a 20 SRR948252_1_cutadapt_STARAligned.sortedByCoord.out.bam /home/agro5431/shared/Zea_may6.AGPv4.33.gff3 > htseq_SRR948252.txt
```

3) Remove adapter sequence

```
#!/bin/bash
#SBATCH --time=3:00:00
#SBATCH --ntasks=8
#SBATCH --mem=16g
#SBATCH --tmp=16g
#SBATCH --mail-type=ALL
#SBATCH --mail-user=cnhirsch@um.edu

echo "Step 0: Setting up directory structure..."
mkdir /home/agro5431/cnhirsch/transcriptome_analysis/
cd /home/agro5431/cnhirsch/transcriptome_analysis/

echo "Step 1: Downloading the file from the SRA..."
module load sratoolkit/2.8.2
fastq-dump --split-files -f Q33 SRR948252

echo "Step 2: Check the quality of the raw reads..."
module load fastqc/0.11.7
fastqc -f fastq SRR948252.1.fastq

echo "Step 3: Remove the adapter sequences..."
module load cutadapt/1.18
cutadapt -a GATCGGAAGAGCACACGTCTGAACTCCAGTCACGATCATCTGATCGCTTCCTTCGTTG -a AATGATACGGCAGCAGACGATCTACATCTTCCCTACACGACGCTTCGTCATC -f fastq -m 30 --quality-base=33 -o SRR948252_1_cutadapt.fastq SRR948252.1.fastq > SRR948252_cutadapt.log

echo "Step 4: Check the quality of the cleaned reads..."
fastqc -f fastq SRR948252_1_cutadapt.fastq

echo "Step 5: Align RNAseq reads to the genome..."
module load star/2.5.3a
STAR --genomeDir /home/agro5431/shared/B73_STAR \
--runThreadN 8 \
--readFilesIn SRR948252_1_cutadapt.fastq \
--outFileNamePrefix SRR948252_1_cutadapt_STAR \
--outSAMtype BAM SortedByCoordinate \
--outSAMunmapped Within \
--outSAMattributes Standard

echo "Step 6: Get transcript abundance..."
module load htseq/0.7.2
htseq-count -f bam -s no -t gene -i ID -m union -a 20 SRR948252_1_cutadapt_STARAligned.sortedByCoord.out.bam /home/agro5431/shared/Zea_may6.AGPv4.33.gff3 > htseq_SRR948252.txt
```

4) Align RNAseq reads to the genome

HISAT2

graph-based alignment of next generation sequencing reads to a population of genomes



HISAT2 is a fast and sensitive alignment program for mapping next-generation sequencing reads (both DNA and RNA) to a population of human genomes (as well as to a single reference genome). Based on an extension of BWT for graphs [Srin et al. 2014], we designed and implemented a graph FM index (GFM), an original approach and its first implementation to the best of our knowledge. In addition to using one global GFM index that represents a population of human genomes, HISAT2 uses a large set of small GFM indexes that collectively cover the whole genome (each index representing a genomic region of 56 Kbp, with 55,000 indexes needed to cover the human population). These small indexes (called local indexes), combined with several alignment strategies, enable rapid and accurate alignment of sequencing reads. This new indexing scheme is called a Hierarchical Graph FM index (HGFMI).



BIOINFORMATICS ORIGINAL PAPER

Vol. 29 no. 1, 2013, pages 15–21
doi:10.1093/bioinformatics/bts635

Sequence analysis

Advance Access publication October 25, 2012

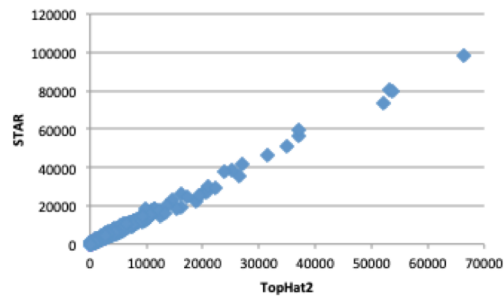
STAR: ultrafast universal RNA-seq aligner

Alexander Dobin¹*, Carrie A. Davis¹, Felix Schlesinger¹, Jorg Drenkow¹, Chris Zaleski¹, Sonali Jha¹, Philippe Batut¹, Mark Chaisson² and Thomas R. Gingeras¹

¹Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and ²Pacific Biosciences, Menlo Park, CA, USA

Associate Editor: Inanc Birol

Comparison of Counts from Different Aligners



4) Align RNAseq reads to the genome

module load star

```
STAR --runMode genomeGenerate \
--genomeDir /home/hirschcl/shared/projects/te-expr/data/databases/star \
--genomeFastaFiles /home/maize/shared/databases/genomes/Zea_mays/PH207/ZmaysPH207_443_v1.0.fa \
--sjdbGTFfile /home/maize/shared/databases/genomes/Zea_mays/PH207/ZmaysPH207_443_v1.1.gene_exons.gtf \
--sjdbOverhang 49 \
--runThreadN 1 \
--limitGenomeGenerateRAM 35773245482
```

4) Align RNAseq reads to the genome

```
#!/bin/bash
#SBATCH --time=3:00:00
#SBATCH --ntasks=8
#SBATCH --mem=16g
#SBATCH --topo=big
#SBATCH --mail-type=ALL
#SBATCH --mail-user=cnhirsch@um.edu

echo "Step 0: Setting up directory:"
mkdir /home/agro5431/cnhirsch/transcriptome
cd /home/agro5431/cnhirsch/transcriptome

echo "Step 1: Downloading the file..."
module load sratoolkit/2.8.2
fastq-dump --split-files -f Q33 SRR948252

echo "Step 2: Check the quality of the raw reads..."
module load fastqc/0.11.7
fastqc -f fastq SRR948252.1.fastq

echo "Step 3: Remove the adapter sequences..."
module load cutadapt/1.18
cutadapt -a GATCGGAAGAGCACACGTCTGAACTCCAGACACGATCATCTGATCGGTCTCTTCTGCTTGA -a AATGATACGGGAGACACGAGATCTACACTCTTCCTACAGACGCTCTCCGATCT -f fastq -m 30 -q quality-base=33 -o SRR948252.1_cutadapt.fastq SRR948252.1.fastq > SRR948252.1_cutadapt.log

echo "Step 4: Check the quality of the cleaned reads..."
fastqc -f fastq SRR948252.1_cutadapt.fastq

echo "Step 5: Align RNAseq reads to the genome..."
module load star/2.5.3a
STAR --genomeDir /home/agro5431/shared/B73_STAR \
--runThreadN 8 \
--readFilesIn SRR948252.1_cutadapt.fastq \
--outFileNamePrefix SRR948252.1_cutadapt_STAR \
--outSAMtype BAM SortedByCoordinate \
--outSAMunmapped Within \
--outSAMattributes Standard

echo "Step 6: Get transcript abundance..."
module load htseq/0.7.2
htseq-count -f bam -s no -t gene -i ID -a union -a 20 SRR948252.1_cutadapt_STARAligned.sortedByCoord.out.bam /home/agro5431/shared/Zea_may
6_429v4_35.gff3 > htseq_SRR948252.txt
```

4) Align RNAseq reads to the genome

% more XXXXXXXX_1_cutadapt_STARLog.final.out

Started job on Apr 01 10:55:10	
Started mapping on Apr 01 10:56:24	
Finished on Apr 01 10:58:56	
Mapping speed, Million of reads per hour	292.59
Number of input reads	12353839
Average input read length	99
UNIQUE READS:	
Uniquely mapped reads number	10636220
Uniquely mapped reads %	86.10%
Average mapped length	99.10
Number of splices: Total	2571479
Number of splices: Annotated (sjdb)	2463999
Number of splices: GT/AG	2530253
Number of splices: GC/AG	34963
Number of splices: AT/AC	1520
Number of splices: Non-canonical	4743
Missmatch rate per base, %	1.55%
Deletion rate per base	0.00%
Deletion average length	1.86
Insertion rate per base	0.00%
Insertion average length	1.38
MULTI-MAPPING READS:	
Number of reads mapped to multiple loci	545476
% of reads mapped to multiple loci	4.42%
Number of reads mapped to too many loci	279225
% of reads mapped to too many loci	2.26%
UNMAPPED READS:	
% of reads unmapped: too many mismatches	0.00%
% of reads unmapped: too short	7.02%
% of reads unmapped: other	0.21%
CHIMERIC READS:	
Number of chimeric reads	0
% of chimeric reads	0.00%

5) Get transcript abundances

5) Get transcript abundances

```
#!/bin/bash
#SBATCH --time=3-00:00
#SBATCH --ntasks=8
#SBATCH --mem=16G
#SBATCH --tmp=16G
#SBATCH --mail-type=ALL
#SBATCH --mail-user=cn

% htseq-count -f bam -s no -t gene -i ID -m union -a 20
XXXXXXXXXX_1_cutadapt_STARAligned.sortedByCoord.out.bam
/home/agro5431/shared/Zea_mays.AGPv4.33.gff3 >
htseq_XXXXXXXXXX.txt

echo "Step 0: Setting
mkdir /home/agro5431/...
cd /home/agro5431/transcriptome_analysis/

echo "Step 1: Downloading the file from the SRA..."
module load sra-toolkit/2.8.2
fastq-dump --split-files -Q 33 SRR948252

echo "Step 2: Check the quality of the raw reads..."
module load fastqc/0.11.7
fastqc -f fastq SRR948252.1.fastq

echo "Step 3: Remove the adapter sequences..."
module load cutadapt/1.18
cutadapt -a GATCGGAAGAGCACACGTCTGAACTCCAGTCACGATGATCTCGTATCCGCTCTTCTGCTTG -a AATGATACGGCAGCAGCAGATCTACATCTTCCCTACACGACGCTCTGCTCAGT
-f fastq -m 30 --quality-base=33 -o SRR948252_1_cutadapt.fastq SRR948252.1.fastq > SRR948252_1_cutadapt.log

echo "Step 4: Check the quality of the cleaned reads..."
fastqc -f fastq SRR948252_1_cutadapt.fastq

echo "Step 5: Align RNAseq reads to the genome..."
module load star/2.5.3a
STAR --genomeDir /home/agro5431/shared/B73_STAR \
--runThreadN 8 \
--readFilesIn SRR948252_1_cutadapt.fastq \
--outFileNamePrefix SRR948252_1_cutadapt_STAR \
--outSAMtype BAM SortedByCoordinate \
--outSAMunmapped Within \
--outSAMattributes Standard

echo "Step 6: Get transcript abundance..."
module load htseq/0.7.2
htseq-count -f bam -s no -t gene -i ID -m union -a 20 SRR948252_1_cutadapt_STARAligned.sortedByCoord.out.bam /home/agro5431/shared/Zea_mays.AGPv4.33.gff3 > htseq_SRR948252.txt
```

6) QC Individual Count Files

- Check that genes have counts
- Check the number of genes printed out is correct (n=39,498)
- use unix 'tail' to look at metrics at bottom of file

```
_no_feature 590886
_ambiguous 132898
_too_low_aQual 0
_not_aligned 745432
_alignment_not_unique 2166383
```

- Look at the annotation of the highest expressed genes


```
% awk '$1 !~ /_/' htseq_XXXXXXXXXX.txt | sort -k 2 -n -r | head -n 3
```

```
% grep <gene_name> /home/agro5431/shared/B73v4.gene_function.txt
```
- Does this make sense?

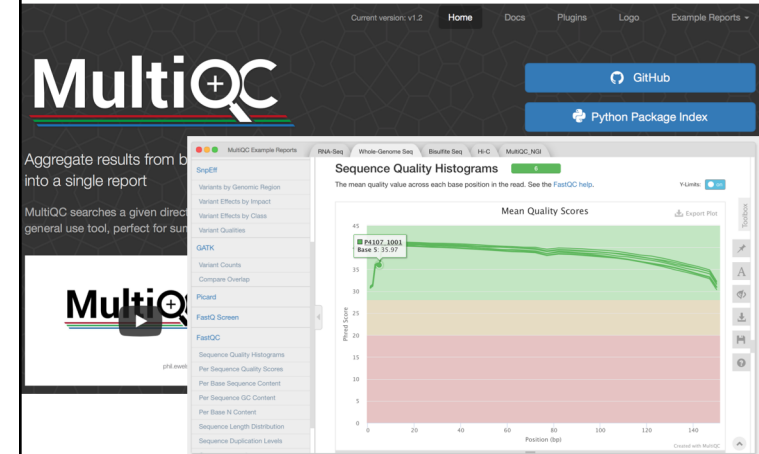
What we are doing today

- 1) Evaluating QC of raw data
- 2) Evaluating QC of read cleaning
- 3) Evaluating QC of read mapping
- 4) Evaluating QC of our matrix
- 5) Determining differentially expressed genes between two tissues

What we are doing today

1) Evaluating QC of raw data

Aggregating Results for QC



What we are doing today

- 1) Evaluating QC of raw data
- 2) Evaluating QC of read cleaning
- 3) Evaluating QC of read mapping

Summarized Statistics from Steps 1-5

The screenshot shows a spreadsheet titled 'summary_statistics'. The spreadsheet has columns labeled A through K. The data is organized into rows, with each row representing a sample. The columns contain various statistics related to sequencing, such as 'SampleName', 'Threat Information', 'Replication', 'Starting Read Number', 'Read Length', 'Reads After Cleaning', 'Mapped Unique', 'Mapped Multiple', 'Unmapped', and 'TruSeq Adapter'. The data is presented in a clear, tabular format, allowing for easy comparison of results across different samples and steps.

Statistics in summary_statistics.xlsx file on GitHub

- ## What we are doing today

Discussion Questions

- What if any samples are problematic?
- What evidence is there for this?
- How many differentially expressed genes did you find?
- Is this more/less than you expected?

[illegible]