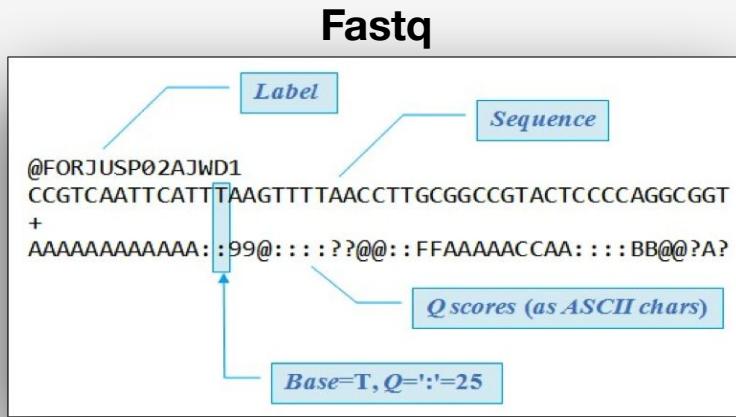


Transcriptomics in Crop Research

Introduction to RNAseq technology

Mary-Ann Blätke
JJ Szymanski

Fastq -> gene expression



Kallisto

-
V

Gene expression

	A	B	C	D
1	genelD	cond1	cond2	cond3
2	ADD3	5.46	5.16	5.14
3	ADGRE1	4.24	4.17	4.21
4	ADGRE5	4.97	5.09	5.03
5	ADGRG6	6.10	6.13	5.98
6	ADGRL2	3.84	4.15	4.26
7	ADH5	3.83	3.75	3.67
8	ADHFE1	6.25	6.11	6.12
9	ADI1	6.23	5.89	5.93
10	ADIPOR1	7.45	7.30	7.26

Our dataset

 
AoBP PLANTS 2020, Vol. 12, No. 5
doi:10.1093/aobpla/plaa041
Advance Access Publication August 19, 2020
Studies

STUDIES

Transcriptome analysis in osmo-primed tomato seeds with enhanced longevity by heat shock treatment

Thiago Barbosa Batista¹, Geysson Javier Fernandez², Tiago Alexandre da Silva¹, Júlio Maia¹ and Edvaldo Aparecido Amaral da Silva^{1*}

¹Department of Plant Production, São Paulo State University (UNESP), Botucatu, São Paulo, Brazil, ²Institute of Biology, Antioquia University, Medellín, Antioquia, Colombia

*Corresponding author's e-mail address: amral.silva@unesp.br

Associate Editor: Gabriela Auge
Form & Function. Chief Editor: Kate McCulloch

Abstract

Seed priming is widely used in commercial seeds and its main function is to accelerate and synchronize seed germination. Undesirably, primed seeds show reduced longevity and treatments like heat shock have been shown to improve longevity in primed seeds. Nonetheless, the effect of heat shock treatment on primed seeds at the mRNA level is not known. Thus, the aim of this work was to investigate the effect of heat shock treatment on the longevity of primed tomato (*Solanum lycopersicum*) seeds at the physiological and transcriptome levels. Tomato seeds were primed and dried (control). Alternatively, primed seeds were subjected to heat shock treatment (38 °C/32 % relative humidity) before drying. Germination, vigor and longevity were evaluated. Transcriptome analysis was performed by RNA sequencing (RNA-seq) from biological samples collected immediately after priming and another samples collected from primed seeds followed by the heat shock treatment. The gene expression was validated by quantitative real time PCR (RT-qPCR). We showed that applying heat shock treatment after priming increased germination speed, enhanced seed longevity and preserved the vigor during storage of primed tomato seeds. Through transcriptome analysis, 368 differentially expressed genes were identified, from which 298 genes were up-regulated and 70 were down-regulated. We showed the increase of mRNA levels of HEAT SHOCK FACTOR-like and HEAT SHOCK PROTEIN-like chaperone genes, suggesting the involvement of the proteins coded by these transcripts in the enhancement of longevity in primed tomato seeds. The heat shock treatment after priming enhances and preserves the vigor of tomato prime seeds during storage. In addition, improves seed longevity through the increase in the expression of transcripts related protection by response to stress.

Keywords: Chaperone molecules; improved longevity; primed seed; seed conservation; seed quality; *Solanum lycopersicum*; storage.

Physiological assays

Seed germination and vigor. Four replications of 50 seeds were germinated in 9 cm Petri dishes with substrate of paper towel moistened with distilled water equivalent to 2.5 times its weight, at 25 °C, under 8 h of light and 16 h in the dark. The length of the primary root, ≥2mm was used as the germination criterion. Data collection was done in different times after sowing; and ended when the germination rate reached 100 % or at 14 days. Seed vigor was determined by the calculation of the time to 50 % of germination (t50) through the analysis of cumulative germination data using the curve fitting module of the Germinator software package (Joosen et al. 2010).

Longevity. We used ageing protocol to assess seed longevity in which the seeds were placed in a support over a saturated solution of NaCl (75 % RH) at 35 °C in glass bottles hermetically sealed. During storage, the water content of *S. lycopersicum* seeds stabilized at $0.10 \pm 0.007 \text{ g H}_2\text{O/g DW}^{-1}$, corresponding to $\pm 9.5 \%$ on wet basis. At different time spans, seeds were imbibed and viability was assessed using the germination assay as described earlier. The different time spans were carried out considering the viability loss behaviour of each treatment group during storage. The viability data were transformed into probit to

libraries were 100 base pair (bp) paired-end sequenced. The data output in fastq file format contained sequence information, including the sequencing quality (Phred quality score). Average Phred scores of ≥20 per position were used for the alignment.

Read alignment and differentially expressed genes. Paired-end reads for mRNA were mapped to the *Solanum lycopersicum* release 39 reference genome using the default parameters of TopHat2 (Kim et al. 2013). Counts for RefSeq genes were obtained using HTSeq (Anders et al. 2015) and DESeq2 (Love et al. 2014) was used to normalize expression counts. The changes in gene expression were considered statistically significant when fold change ≥2 and P-values ≤ 0.05. The RNAseq data was deposited in NCBI (BioProject PRJNA562700: <https://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP220280>).

The analysis of principal components was made using all the genes expressed on the RNA seq data. The normalized count per gene was used and transformed to Z-score. This matrix was used to perform the PCA. For plotting the PCA results, we used the principal component one and two. The heatmap was generated using the normalized counts of the differentially expressed genes. Then we transformed it to z-score and plotted it using the package pheatmap of R.

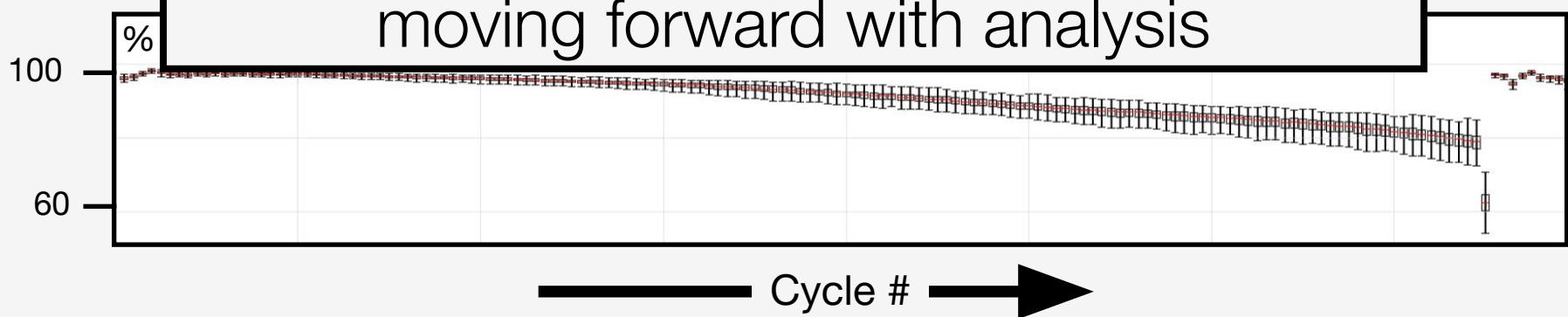
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7494243/>

Read quality

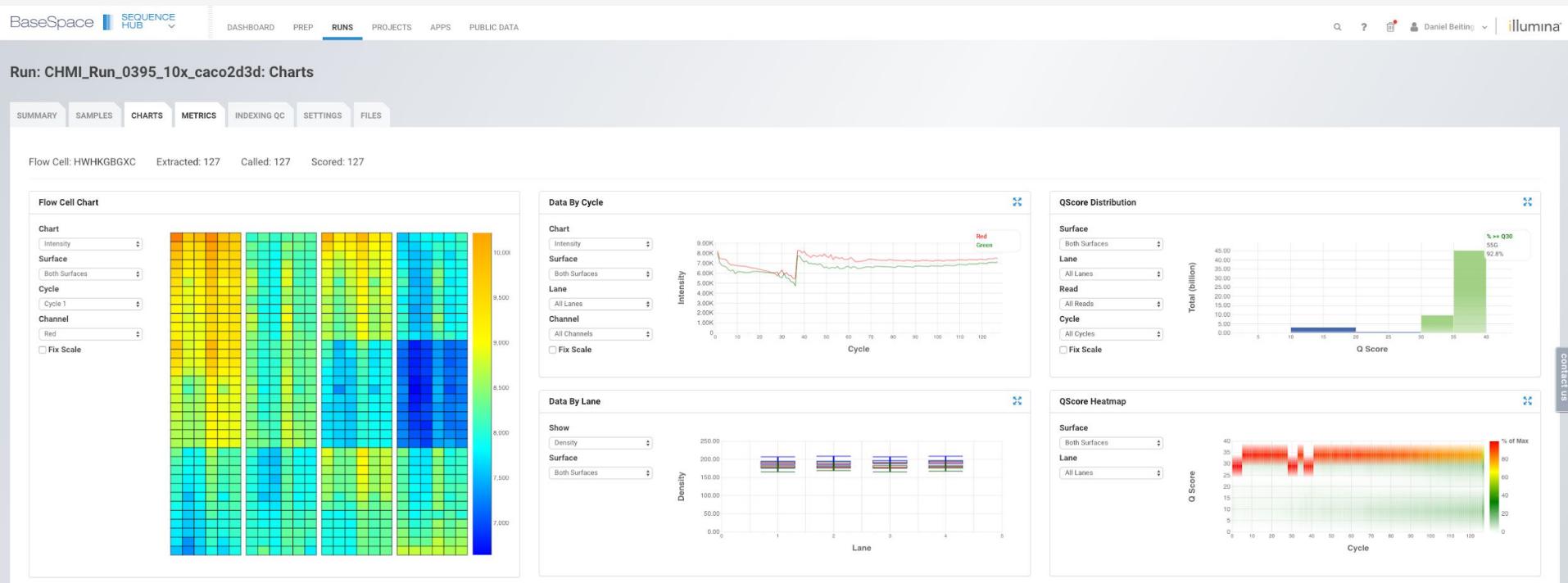
Fastq

		ASCII characters																												
		-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	@	A	B	C	D	E	F	G	H	I
@FORJUSP02AJWD1		1	1	1	1	1	1	1	1	2	2	2	2	2	2	2	2	2	2	3	3	3	3	3	3	3	3	4		
CCGTCATTACATTAAAGTTAACCTTGCAGCGTACTCCCCAGGGT		2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	0

Assess the quality of your raw data before moving forward with analysis



Quality assessment via Illumina's BaseSpace



Quality assessment of fastq files using fastqc



Babraham Bioinformatics



About | People | Services | Projects | Training | Publications

FastQC

Function	A quality control tool for high throughput sequence data.
Language	Java
Requirements	A suitable Java Runtime Environment The Picard BAM/SAM Libraries (included in download)
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

[Download Now](#)

<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Quality assessment of fastq files using fastqc

The screenshot shows a web browser window with the Bhabham Bioinformatics website on the left and a macOS System Report window on the right.

Babraham Bioinformatics

FastQC

Function	A quality control tool for high throughput sequencing data.
Language	Java
Requirements	A suitable Java Runtime Environment . The Picard BAM/SAM Libraries (included in the distribution).
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

[Download](#)

macOS Monterey
Version 12.1

MacBook Pro (13-inch, 2020, Four Thunderbolt 3 ports)

Processor 2.3 GHz Quad-Core Intel Core i7

Memory 32 GB 3733 MHz LPDDR4X

Graphics Intel Iris Plus Graphics 1536 MB

Serial Number [REDACTED]

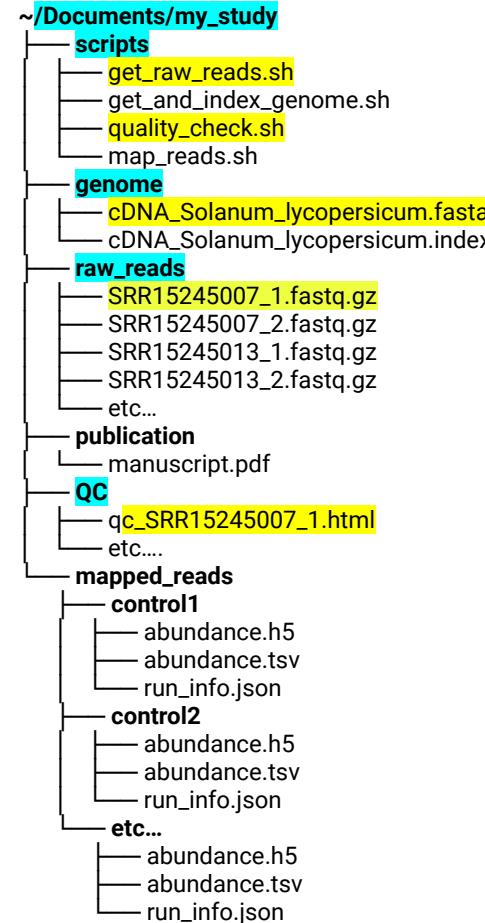
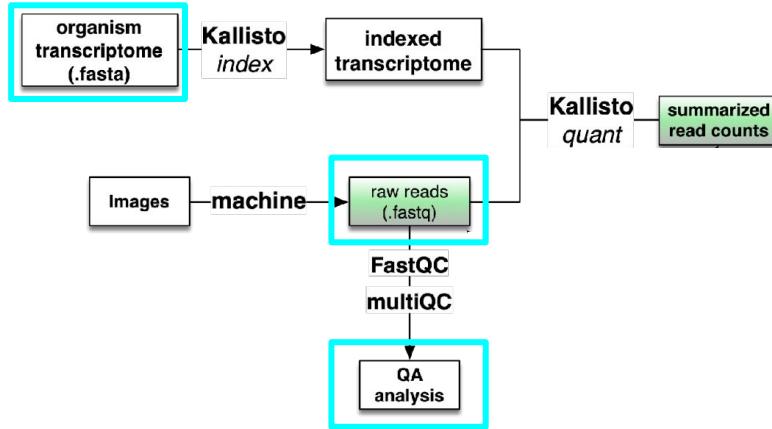
[System Report...](#) [Software Update...](#)

™ and © 1983-2021 Apple Inc. All Rights Reserved. License and Warranty

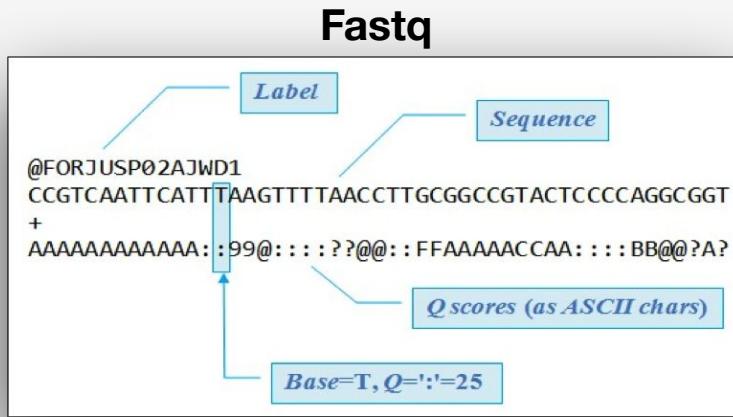
fastqc demonstration

```
fastqc -o ~/my_path/QC/ -t 8 ~/my_path/raw_reads*.fastq.gz
```

RNAseq analysis pipeline



Fastq -> gene expression



Kallisto
->

Gene expression

	A	B	C	D
1	genelD	cond1	cond2	cond3
2	ADD3	5.46	5.16	5.14
3	ADGRE1	4.24	4.17	4.21
4	ADGRE5	4.97	5.09	5.03
5	ADGRG6	6.10	6.13	5.98
6	ADGRL2	3.84	4.15	4.26
7	ADH5	3.83	3.75	3.67
8	ADHFE1	6.25	6.11	6.12
9	ADI1	6.23	5.89	5.93
10	ADIPOR1	7.45	7.30	7.26

Ingredients for an alignment

our data
(or your data)

.fastq

reference
sequences

.fasta

Ingredients for an alignment

our data

```
1 @SRR10056916.5114340 5114340 length=101
2 GATTAACATATTTCATTTGAACTTGTAGCGTTCTACACATTGCTTACCCCTCCCTCGATAATATCCACCGAACATCCGTGG
3
4 @SRR10056916.5114341 5114341 length=101
5 CGCCGTTCCCGAGAGGGTGTTCAGAGCTGGAGTTATCCAAGAGATAACGCCGAAGATTTTTGAATCTACTCAAGAAAAGAAGAAAAGAAGTGTGG
6
7 @SRR10056916.5114342 5114342 length=101
8 GTGAAAGTACTACGCCCTGGTGGCCGGAGAGGTAATGGTGTGGATATGGTGAATCAGCAGTGTTCAGCCCCATCC
9
10 @SRR10056916.5114343 5114343 length=101
11 CTAGGAGATTTTGATCTCATCAGAGTTATACGTAATAGAAAAGTAGAGATGACCTAAGAGTAAGAATCTAATATTGAATATTGAATGAGGG
12
13 @SRR10056916.5114344 5114344 length=100
14 AGACAGACGTGTTTGCCTGCAGAACTGCAACGGGATAGTGAGGGGCCACCGAGGAGGTTGATGGATGCTATGGCTAGGCAGACTCCACTGGG
15
16 @SRR10056916.5114345 5114345 length=99
17 AGTATGTTAACCTATCGCGAACATTCAAACAAAAGGGAAAGAGGTTAGCAAGGCTCTAAGATCCAGAGGCTGGAACACCATTGACTCTC
18
19 @SRR10056916.5114346 5114346 length=100
20 ACTAGAACATCATCACCCACTTACATCTCAGCAAAACCATCATGTAACAGTGATGGTGCCTCTGGAGAGCGTGGATCAAATGCACGAACAGTG
21
22 @SRR10056916.5114347 5114347 length=101
23 GCTTAAGAGGTAGAAAGAGCAACGGTCAGATGCTGGCTTAAATAGAGATTTTGAGAAGAACCAACAACATCAAGAATTATGGATTTGGCTCCGT
24
25 @SRR10056916.5114348 5114348 length=101
26 TGAGTGGCTCTGGATTCTGCCACACTCTCGATTCCCCGGTTCATCTGAACATATAATGTGGAGAAATGCCATGAGGAAGGAGGCCAGTATT
```

.fastq

reference

```
ACGTTCTGGAGTGGCTGCAGCTACCCATGATGCGTAACAGGCTGGCCATCCAAG
CCATGCAGGATCACTCAGGATTTCAGTTCACCTCTATTCCAAGCATTACCTCAA
AGGACCCAGCAGCTACACCCCTACAGGCTTCAAGGCCCTCATAGTCATGCTCTCC
CATTTACCCCTACCCATCCTGATCGGTATGCCCTAGCCTGACCCCTTAGATAAGCAA
TGAGGTAGGAAGAACAAACCCCTTGGCTTCCCTGGCATGTTGGAGAAAGTGCCTGCCTGG
TCCGAGCCGCCCTGGTCTGAAGCAGGTGCTCTGCTTACCTTGCTTAGGCTGCTGCA
GAAGCACCTGCCGTGCACTCAGCACCTCTTGCTAGAGCCCTCATCACCTCAGG
CTGTCACCATGGGCCAGGAACCAACAGCACTGGTTACTGCTGGGGTAAACT
AACTCAGTGGAAATGGGCTGTTACTTGGCTGCTCAACTCATAAAGTTGGCTGATT
TGAAAAAAAGCTCATCAAATAAAAGGCTATGTTGCTGGCTGGTCCC
>ENST00000397163.8 cdna chromosome:GRCh38:15:42359501:42412317:1 gene:ENSG0000092529.25
gene_biotype:protein_coding transcript_transcript_biotype:protein_coding gene_symbol:CAPN3 description:calpain 3
[Source:HGN Symbol;Acc:HGN;C:1480]
ACTCTTTCTCTCCCTCTGGCATGCTGCTGGAGAGCCCAAGTCACAT
TGCTTCAGAAATCTTAACTGACTCATTTCTCAGGAGAACTTATGGCTTCAGAACTCACAGC
TCGTTTTAAAGATGGACATAACCTGAGACCTCTGATGGCTTCACATTGAACTG
GATGTGGACACTTTCTCAGATGACAGAAGATTACTCCAACTTCCCTTGCAGTTGCTT
CCTTCCTTGAAGGTTAGCTGTATCTTATTCTTAAAGCTTTCTCAGGAAAGCCAC
TTGCCATGCCGACCGTCTTGGCTCCAAAGGCAAGGGCTGAGGCCGGT
CCCCAGGGCCAGTCTTACCCGGCCAGAGCAAGGGCACTGAGGCTGGGGTGGAAACC
CAAGTGGCATATTACGGCATCATGCCAACTTTCTATTATCGGAGTAAAGAGA
AGACATTGGAGCAACCTCACAAGAAATGCTAGAAAAGAAGTTCTTATGTGGACCTG
AGTCCCCACGGATGAGACCTCTCTTTTATAGCCAGAAGTTCCCATCAGTTGCT
GGAAAGACCTCCGGAAATTGGAGAATCCCCGATTATCATGATGGAGCAACAGAA
CTGACATCTGTCAGGGAGGCTAGGGGACTGCTGGTTCTGGAGGCCATTGCTGCTGA
CCCTGAACCAAGCACCTTTCGAGTCATACCCATGATCAAAGTTCATGAAACAT
```

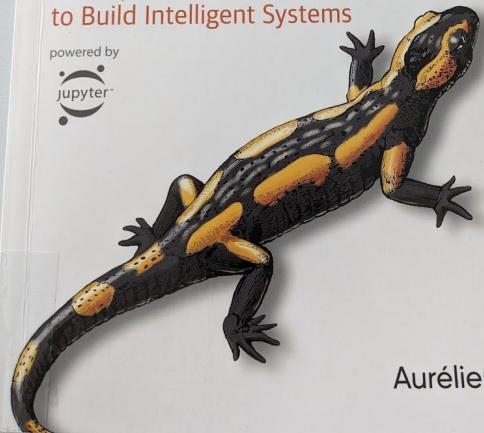
.fasta

O'REILLY®

Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow

Concepts, Tools, and Techniques
to Build Intelligent Systems

powered by
 jupyter™



Aurélien Géron

2nd Edition
Updated for
TensorFlow 2

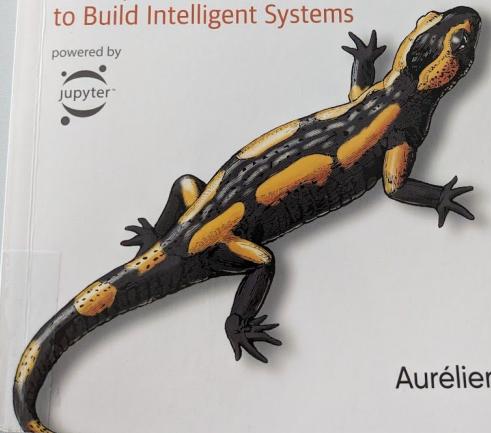
“tensor arrays”

O'REILLY®

Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow

Concepts, Tools, and Techniques
to Build Intelligent Systems

powered by
 jupyter™



Aurélien Géron

2nd Edition
Updated for
TensorFlow 2

T

- t-Distributed Stochastic Neighbor Embedding (t-SNE), 233
- tail-heavy histograms, 51
- Talos library, 322
- target model, 639
- TD error, 630
- TD target, 630
- temperature
 - in Boltzmann machines, 775
 - in text generation, 531
- Temporal Difference Learning (TD Learning), 629
- tensor arrays, 383, 786
- TensorBoard, 317

Ingredients for an alignment

our data
(or your data)

.fastq

reference
sequences
index

.fasta

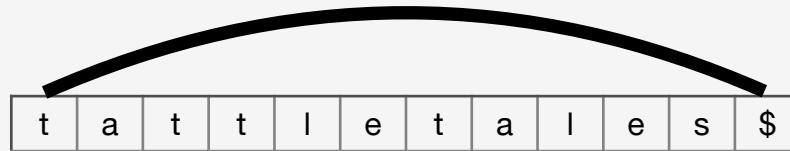
Building an index

- Use Kallisto to build index from reference fasta
 - `kallisto index -i inputFastaName.index inputFastaName.fa`
 - `kallisto index -i Homo_sapiens.GRCh38.cdna.all.index Homo_sapiens.GRCh38.cdna.all.fa`
- Only need to build an index once
- Be careful with file names!
 - Long and meaningful file names are better than short ambiguous ones.
 - No spaces. ‘`my_file_name`’, or ‘`myFileName`’, or ‘`my.file.name`’....just not, ‘`my file name`’

Let’s construct an index from our reference fasta file

Burrows-Wheeler Transform (BWT)

‘tattletales’



Burrows-Wheeler Transform (BWT)

‘tattletales’

12	\$	t	a	t	t	l	e	t	a	l	e	s
8	a	l	e	s	\$	t	a	t	t	l	e	t
2	a	t	t	l	e	t	a	l	e	s	\$	t
10	e	s	\$	t	a	t	t	l	e	t	a	l
6	e	t	a	l	e	s	\$	t	a	t	t	l
9	l	e	s	\$	t	a	t	t	l	e	t	a
5	l	e	t	a	l	e	s	\$	t	a	t	t
11	s	\$	t	a	t	t	l	e	t	a	l	e
7	t	a	l	e	s	\$	t	a	t	t	l	e
1	t	a	t	t	l	e	t	a	l	e	s	\$
4	t	l	e	t	a	l	e	s	\$	t	a	t
3	t	t	l	e	t	a	l	e	s	\$	t	a

Burrows-Wheeler Transform (BWT)

‘tattletales’

		Firs										Las	
		t	\$	t	a	t	t	l	e	t	a	l	e
12		\$	t	a	t	t	l	e	t	a	l	e	t
8		a	l	e	s	\$	t	a	t	t	l	e	t
2		a	t	t	l	e	t	a	l	e	s	\$	t
10		e	s	\$	t	a	t	t	l	e	t	a	t
6		e	t	a	l	e	s	\$	t	a	t	t	t
9		l	e	s	\$	t	a	t	t	l	e	t	t
5		l	e	t	a	l	e	s	\$	t	a	t	t
11		s	\$	t	a	t	t	l	e	t	a	l	e
7		t	a	l	e	s	\$	t	a	t	t	l	e
1		t	a	t	t	l	e	t	a	l	e	s	\$
4		t	l	e	t	a	l	e	s	\$	t	a	t
3		t	t	l	e	t	a	l	e	s	\$	t	a

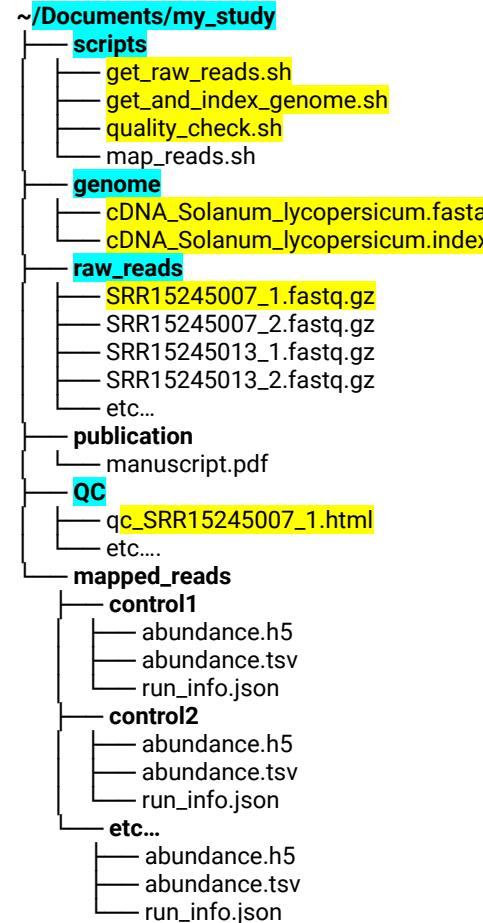
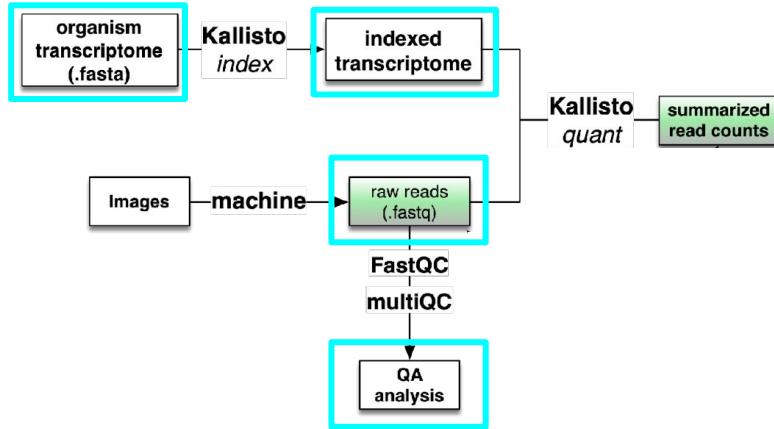
A red vertical rectangle highlights the last column of the grid, which contains the characters 't', 'a', 's', 't', 't', 'l', 'i', 'a', 't', 'e', 't', 'a', '\$', and 't'. To the right of the grid, the word 'tattletales' is shown above a downward-pointing arrow labeled 'BWT'. Below the arrow, the transformed string 'stllatee\$ta' is displayed.

Burrows-Wheeler Transform (BWT)

‘tattletales’

		Firs												
		t	\$	t	a	t	t	l	e	t	a	l	e	s
12		t	\$	t	a	t	t	l	e	t	a	l	e	s
8		a	l	e	s	\$	t	a	t	t	l	e		t₄
2		a	t	t	l	e	t	a	l	e	s	\$		t₁
10		e	s	\$	t	a	t	t	l	e	t	a	l	
6		e	t	a	l	e	s	\$	t	a	t	t	l	
9		l	e	s	\$	t	a	t	t	l	e	t	a	
5		l	e	t	a	l	e	s	\$	t	a	t		t₃
11		s	\$	t	a	t	t	l	e	t	a	l	e	
7		t₄	a	l	e	s	\$	t	a	t	t	l	e	
1		t₁	a	t	t	l	e	t	a	l	e	s	\$	
4		t₃	l	e	t	a	l	e	s	\$	t	a		t₂
3		t₂	t	l	e	t	a	l	e	s	\$	t	a	

RNAseq analysis pipeline



Application of BWT in genomics

BIOINFORMATICS

ORIGINAL PAPER

Vol. 25 no. 14 2009, pages 1754–1760
doi:10.1093/bioinformatics/btp324

Sequence analysis

Fast and accurate short read alignment with Burrows–Wheeler transform

Heng Li and Richard Durbin*

Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cam

Received on February 20, 2009; revised on May 6, 2009; accepted on May 12, 2009

Advance Access publication May 18, 2009

Associate Editor: John Quackenbush

~23,000 citations

Open Access

Software

Ultrafast and memory-efficient alignment of short DNA sequences to the human genome

Ben Langmead, Cole Trapnell, Mihai Pop and Steven L Salzberg

Address: Center for Bioinformatics and Computational Biology, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA.

Correspondence: Ben Langmead. Email: langmead@cs.umd.edu

Published: 4 March 2009

Genome Biology 2009, 10:R25 (doi:10.1186/gb-2009-10-3-r25)

Received: 21 October 2008

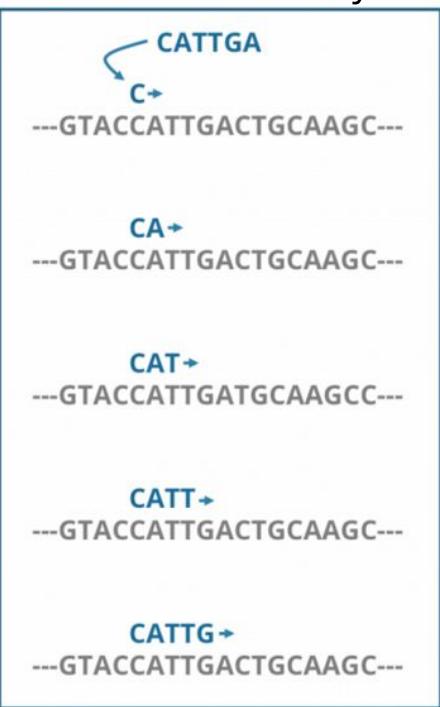
Revised: 19 December 2008

Accepted: 4 March 2009

~16,000 citations

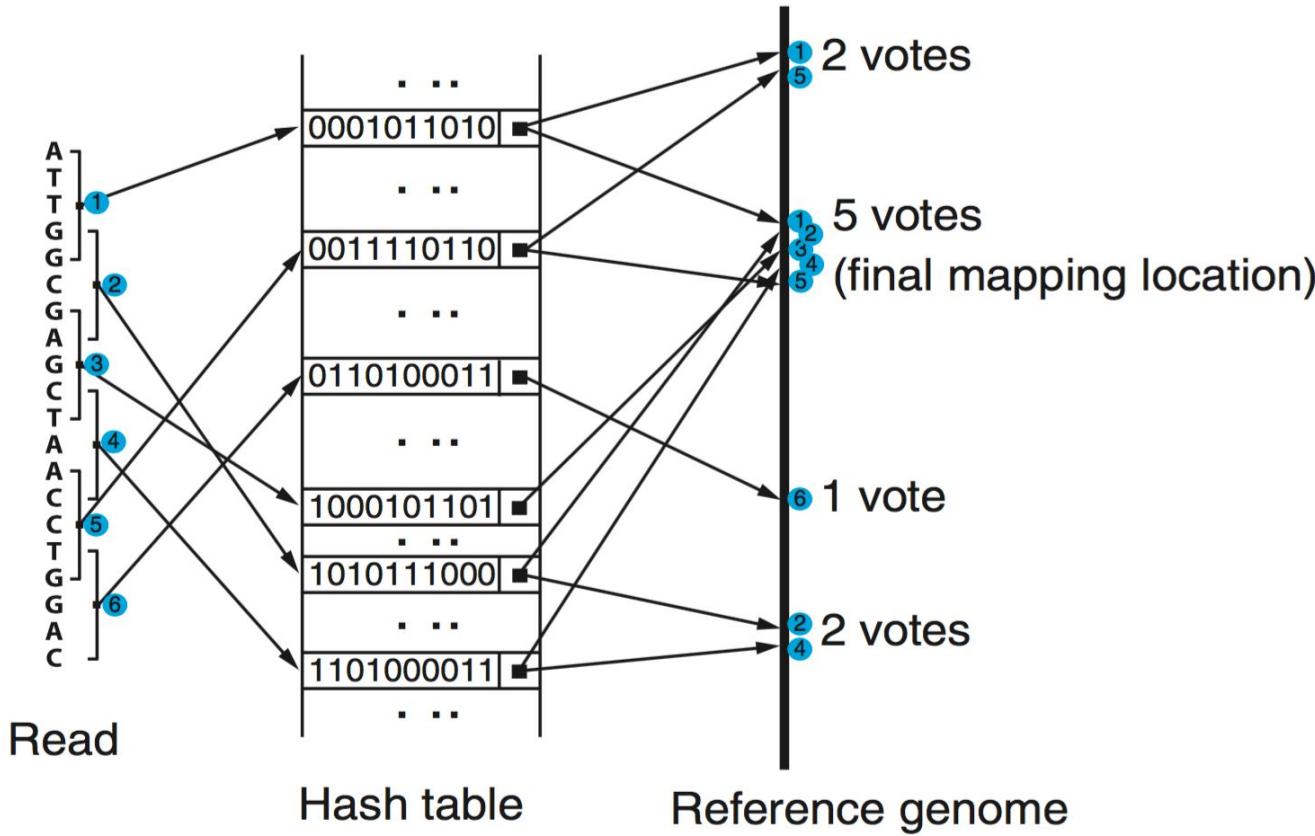
‘Seed and Extend’ alignments

extension only



brute force, but slow

R Subread aligner



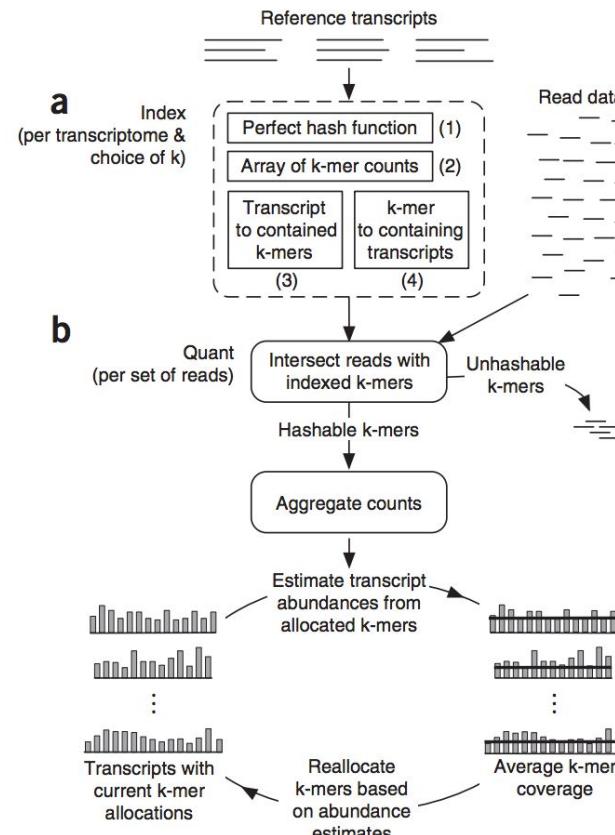
Reads are hard to align, but that's a function of their length. Why not approach the problem by shredding reads into Kmers - enter “Sailfish”

Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms

Rob Patro¹, Stephen M Mount^{2,3} & Carl Kingsford¹

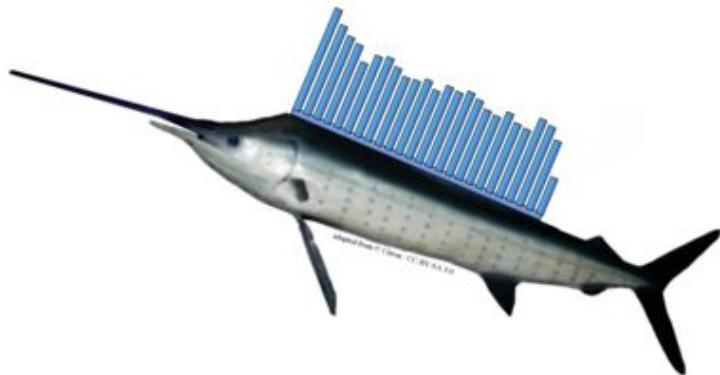
Nature Biotechnology **32**, 462–464 (2014)

Sailfish replaces *approximate* alignment
of (error prone) *reads* with exact
alignment of short *k-mers*



Sailfish

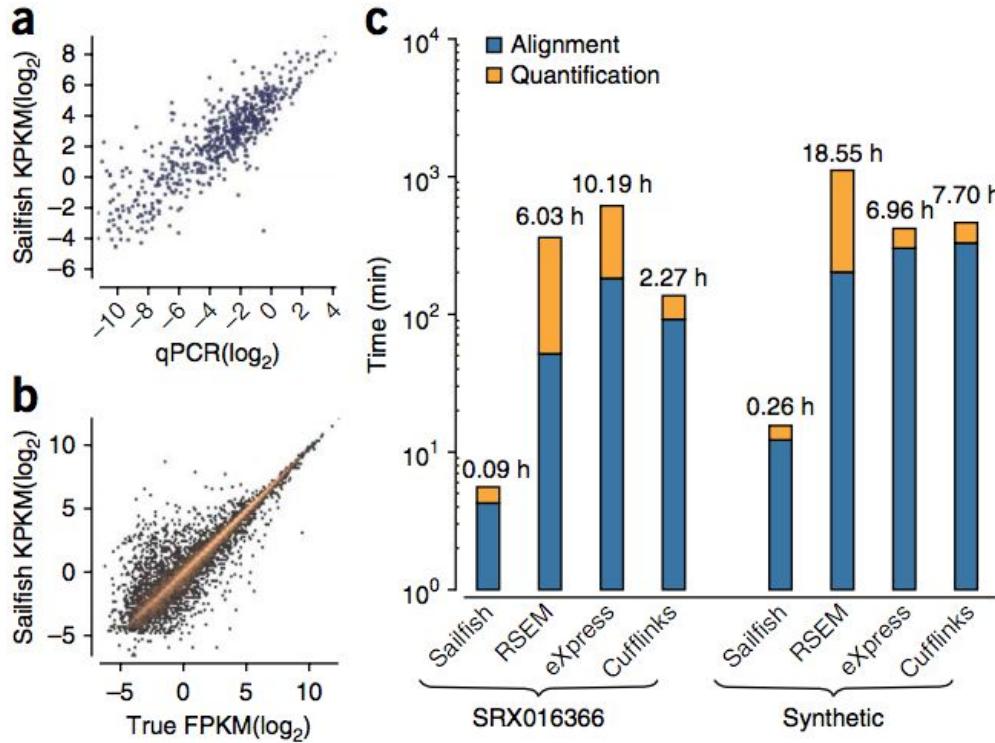
'lightweight read alignment'



“By not requiring read mapping, **Sailfish avoids parameters specifying**, for example, the number of mismatches to tolerate, total allowable quality of mismatched bases, gap open and extension penalties, whether and how much to trim reads, number and quality of alignments to report from the aligner and pass into the estimation procedure.”

Sailfish

'lightweight read alignment'



What if the idea of alignment was altogether abandoned enter “Kallisto”

BRIEF COMMUNICATIONS

nature
biotechnology

Nature Biotechnology **34**, 525–527 (2016)

Near-optimal probabilistic RNA-seq quantification

Nicolas L Bray¹, Harold Pimentel², Páll Melsted³
& Lior Pachter^{2,4,5}

We present kallisto, an RNA-seq quantification program that is two orders of magnitude faster than previous approaches and achieves similar accuracy. Kallisto pseudoaligns reads to a reference, producing a list of transcripts that are compatible with each read while avoiding alignment of individual bases. We use kallisto to analyze 30 million unaligned paired-end RNA-seq reads in <10 min on a standard laptop computer. This removes a major computational bottleneck in RNA-seq analysis.

The first two steps in typical transcript-level RNA-seq processing

this information, we develop a method based on pseudoalignment of reads and fragments, which focuses only on identifying the transcripts from which the reads could have originated and does not try to pinpoint exactly how the sequences of the reads and transcripts align.

A pseudoalignment of a read to a set of transcripts, T , is a subset, $S \subseteq T$, without specific coordinates mapping each base in the read to specific positions in each of the transcripts in S . Accurate pseudoalignments of reads to a transcriptome can be obtained using fast hashing of k -mers together with the transcriptome de Bruijn graph (T-DBG). de Bruijn graphs have been crucial for DNA and RNA assembly⁸, where they are usually constructed from reads. Kallisto uses a T-DBG, which is a de Bruijn graph constructed from k -mers present in the transcriptome (Fig. 1a), and a path covering of the graph, a set of paths whose union covers all edges of the graph, where the paths correspond to transcripts (Fig. 1b). This path covering of a T-DBG induces multi-sets on the vertices, called k -compatibility classes. A compatibility class can be associated to an error-free read by representing it as a path in the graph and defining the k -compatibility class for each vertex as the union of the sets of the k -compatibility classes that contain the path. The kallisto algorithm then finds the transcriptome with the fewest transcripts that are compatible with all the reads.



Quantitative Biology > Genomics

Sailfish: Alignment-free Isoform Quantification from RNA-seq Reads using Lightweight Algorithms

Rob Patro (1), Stephen M. Mount (2), Carl Kingsford (1) ((1) Lane Center for Computational Biology, School of Computer Science, Carnegie Mellon University, (2) Department of Cell Biology and Molecular Genetics and Center for Bioinformatics and Computational Biology,

University of Maryland)

(Submitted on 16 Aug 2013)

RNA-seq has rapidly become the de facto technique to measure gene expression. However, the time required for analysis has not kept up with the pace of data generation. Here we introduce Sailfish, a novel computational method for quantifying the abundance of previously annotated RNA isoforms from RNA-seq data. Sailfish entirely avoids mapping reads, which is a time-consuming step in all current methods. Sailfish provides



Cornell University
Library

Current browse context:

q-bio.GN

< prev | next >

new | recent | 1308

Change to browse by:

cs

cs.CE

q-bio

We gratefully acknowledge support from
the Simons Foundation
and University of Pennsylvania



Quantitative Biology > Quantitative Methods

Near-optimal RNA-Seq quantification

Michael D. Bray, Harald T. Irmentent, Pál Melsted, Lior Pachter

(Submitted on 11 May 2015 (v1), last revised 15 May 2015 (this version, v2))

We present a novel approach to RNA-Seq quantification that is near optimal in speed and accuracy. Software implementing the approach, called kallisto, can be used to analyze 30 million unaligned paired-end RNA-Seq reads in less than 5 minutes on a standard laptop computer while providing results as accurate as those of the best existing tools. This removes a major computational bottleneck in RNA-Seq analysis.

Comments: - Added some results (paralog analysis, allele specific expression analysis, alignment comparison, accuracy analysis with TPMs) - Switched bootstrap analysis to human sample from SEQC-MAQCIII - Provided link to a snakefile that allows for reproducibility of all results and figures in the paper

Subjects: Quantitative Methods (q-bio.QM); Computational Engineering, Finance, and Science (cs.CE); Data Structures and Algorithms (cs.DS); Genomics (q-bio.GN)

Cite as: arXiv:1505.02710 [q-bio.QM]
(or arXiv:1505.02710v2 [q-bio.QM] for this version)

Submission history

From: Lior Pachter [view email]

[v1] Mon, 11 May 2015 17:42:04 GMT (3410kb)

[v2] Fri, 15 May 2015 17:12:58 GMT (4940kb)

Download:

- PDF only

(license)

Current browse context:

q-bio.QM

< prev | next >

new | recent | 1505

Change to browse by:

cs

cs.CE

cs.DS

q-bio

q-bio.GN

References & Citations

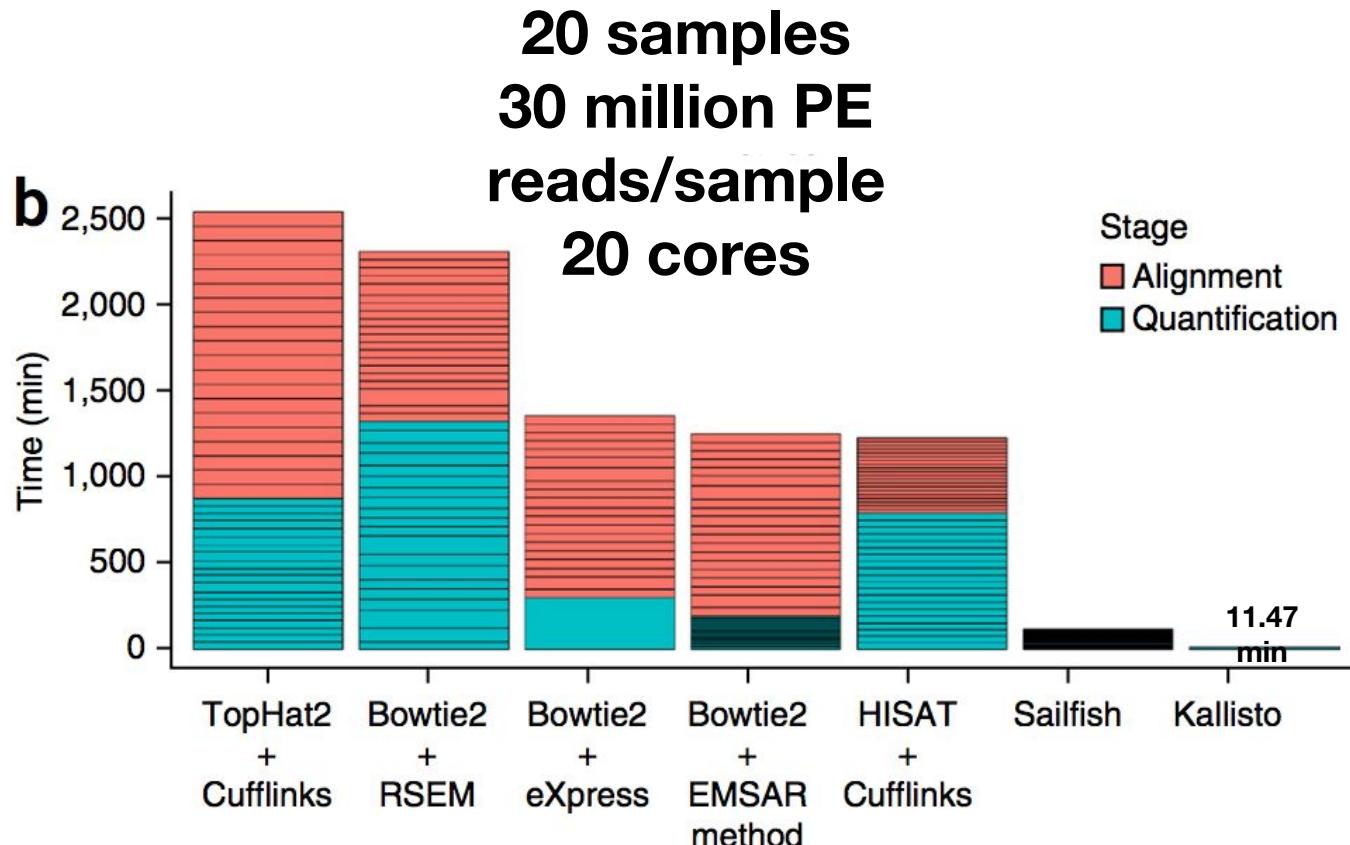
- NASA ADS

Bookmark

(what is this?)



Kallisto - it's fast



Mapping single-end data

Calling the program and function

Name and path of index

Name of output folder

Number of threads to use

Single vs paired-end
and fragment size

Name of fastq file to be mapped

Code continues on next line

```
kallisto quant -i Homo_sapiens.GRCh38.cdna.all.index -o test -t 8 --single -l 250 -s 30 SRR8668755_1M_subsample.fastq.gz
```

Mapping paired-end data

Calling the program and function

Name and path of index

Name of output folder

Number of threads to use

Name of read1 fastq file

Name of read2 fastq file

Code continues on next line

```
• kallisto quant |
  -i Homo_sapiens.GRCh38.cdna.all.index |
  -o test |
  -t 8 |
  sample1_read1.fq.gz |
  sample1_read2.fq.gz |
  &> test.log
```

Mapping single-end data

Calling the program and function

Name and path of index

Name of output folder

Number of threads to use

Single vs paired-end
and fragment size

Name of fastq file to be mapped

Code continues on next line

```
kallisto quant |  
-i Homo_sapiens.GRCh38.cdna.all.index |  
-o test |  
-t 8 |  
--single -l 250 -s 30 |  
SRR8668755_.fastq.gz  
&> test.log
```

Kallisto demonstration

```
kallisto quant |  
-i Homo_sapiens.GRCh38.cdna.all.index |  
-o test |  
-t 8 |  
sample1_read1.fq.gz |  
sample1_read2.fq.gz |  
&> test.log
```