基生研ゲノムインフォマティックス・トレーニングコース RNA-seq入門 - NGSの基礎からde novo解析まで 実践編:RNA-seq解析パイプライン 2017春

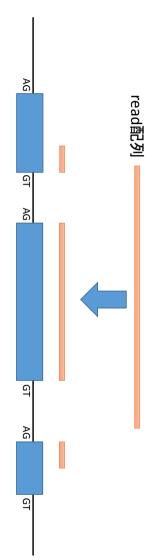
RNA-seq解析バイフレイン 2017.03.09-2017.03.10

KNA-Sedパイプライソゲノスベースの解析法

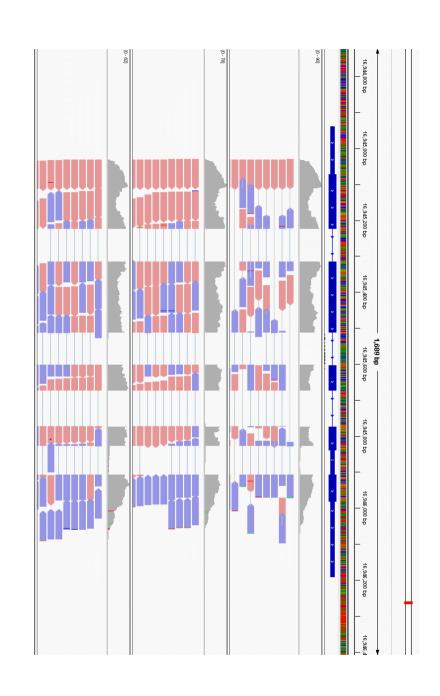
基礎生物学研究所生物機能解析センター山口勝司

genomeをレファレンスとする場合

レファレンスがゲノム配列の場合 イントロン配列のスプライシングを考慮した アライメントを行う必要がある。 TopHatを用いる 他 Blat, SpliceMap, MapSplice, GSMAP, QPALMA



実際こんな感じにアラインされる



本下フーニングコースたの流ち

reads genome Preprocessing し
Tophat (スプライスサイトを考慮したアライメントツール)

アライメント

gene model (GTFファイル)

Cufflinks (カウント、DE gene等の検出等)

DE gene等のリスト

Tophat



Hat 2.1.1 release 2/23/2016

- TopHat source code moved to GitHub 3/31/2015 pHat is now available as a public GitHub repository when the control of the cont



indelを考慮したアライメント Bowtie2/二対応 TopHat2になりalignerとして

が可能になった

Kim et al. Genome Biology 2013, 14:R36 http://genomebiology.com/2013/14/4/R36



Open Access

METHOD

fusions the presence of insertions, deletions and gene TopHat2: accurate alignment of transcriptomes in

Daehwan Kim^{1,2,3*}, Geo Pertea³, Cole Trapnell^{5,6}, Harold Pimentel⁷, Ryan Kelley⁸ and Steven L Salzberg^{3,4}

mapping to known transcripts, producing sensitive and accurate alignments, even for highly repetitive genomes or in the presence of pseudogenes. TopHat2 is available at http://ccb.jhu.edu/software/tophat. can occur after genomic translocations. TopHat2 combines the ability to identify novel splice sites with direct reference genome. In addition to de novo spliced alignment, TopHat2 can align reads across fusion breaks, which produced by the latest sequencing technologies, while allowing for variable-length indels with respect to the TopHat is a popular spliced aligner for RNA-sequence (RNA-seq) experiments. In this paper, we describe TopHat2 which incorporates many significant enhancements to TopHat. TopHat2 can align reads of various lengths

for RNA-S



TopHat is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to m and then analyzes the mapping results to identify splice junctions between exons. omes using the ultra high-throughput short re

TopHat is a collaborative effort among Daehwan Kim and Steven Salzberg in the Center for Computational Biology at Johns Hopkins University, and Cole Trapnell in the Genome Sciences Department at the University of Washington. TopHat was originally developed by Cole Trapnell at the Center for Bioinformatics and Computational Bio at the University of Maryland, College Park.

TopHat 2.1.1 release 2/23/2016

- Please note that TopHat has entered a low maintenance, low support stage as it is now largely superseded by **HISAT2** which provides the same core functionality (i.e. spliced alignment of RNA-Seq reads), in a more accurate and **much more efficient** way.

 Version 2.1.1 is a maintenance release which includes the following changes, some of them thanks to GitHub contributors:

 TopHat can be now built on more recent Linux distributions with newer GNU C++ (5.x), as the included SeqAn library was finally upgraded to
- improved the detection of linker options for the Boost::Thread library which prevented the TopHat build from source on some systems.
- incorporated Luca Venturini's code to support large bowtie2 indexes (.bt2l).

 bam2fastx usage message (-h/--help) was updated in order to better document the functions of this program which can be used as a standalone utility for converting reads from BAM/SAM to FASTQ/FASTA; the -v/--version option was also added to this utility for easier integration in other pipelines.

TopHat 2.1.0 release 6/29/2015

TopHat-Fusion algorithm improves Kelley at Illur ements for more sensitive and accurate discovery of fusions, thanks to contributions from Gordon Bean and

とりあえず使って みる

Tools Users Google Group. Please use tophat.cufflinks@gmail.com for private communications only. Please

Getting startedで

FAQ Protocol

News and updates

Site Map

OSI certified

- This release implements a new algorithm for counting fusion-supporting read pairs that reduces the number of false-positive potential fusions. This algorithm computes the inner distance between read pairs by first converting the pair positions to transcript coordinates using the transcript information in refleene.txt and ensGene.txt. Pairs with small inner distance (suggesting the pair could come from a plausible pair-end insert) are counted as supporting evidence for the fusion. The default threshold for the inner distance is 250 base pairs; this parameter can be set using the —fusion-pair-dist <int> flag.
- fixed a few issues with GFF parsing of some annotation files
 fixed a runtime-error when using --no-discordant option.
 Several fixes/improvements thanks to contributors on GitHub:
 new --max-num-fusions option allowing the user to specify
- sting lower limit for --fusio
- fixed a few typos, cleaning up python code etc.

TopHat source code moved to GitHub 3/31/2015

TopHat is now available as a pusubmit patches (pull requests). to submit bug reports (issues)

TopHat 2.0.14 release 3/24/2015 Version 2.0.14 is a maintenance relea

- Version 2.0.14 is a maintenance release with the following changes:

 pipeline speed improvements thanks to contributions from Véronique Legrand and Michael Press

 added support for xz compressed read files (thanks to a patch submitted by Ashton Trey Belew)

quantitation for RNA-Seq
Bowtie: Ultrafast short read alignment
TopHat-Fusion: An algorithm for
Discovery of Novel Fusion Transcripts
CummeRbund: Visualization of RNA-Related Tools Cufflinks: Isofo Linux x86_64 binary
Mac OS X x86_64 binary

Releases

TopHat contributors directly. do not email technical que

tions to

Getting started

l quick-start ad and extract the latest Bo

simply

Hat version already installed on your system of copying the new programs in a directory directory in your shell's PATH. For example, ider your home directory:

Installing a pre-compiled binary release
In order to make it easy to install Tophat we provide a few binary packar
In order to make it easy to install Tophat we provide a few binary packar
Inophat themselves, which requires a certain development environment at
download the appropriate one for your patform, unpack it, and make use variable (or create a symbolic link to the included tophat2 script somewh
Moret if you want to be able to install and run this new version without it
make sure you unpack the new version into a different directory from the
lin your PATH just create a symbolic link from the sephant2 wrapper scrip
assuming the ~/bin directory is in your PATH and you unpack tophat-2.7 y packages to save users from the occar mment and the Boost libraries installs make sure the reginax binaries are in a d somewhere in your PATH, some seek of without overwriting a pervious Fopblaw) without overwriting a pervious Fopblaw of the open seek of the pervious for path of the open seek of the pervious for path of the open seek of the pervious for seek of the open seek of the pervious for seek of the open seek of the pervious for the open seek of the open seek of the pervious for the open seek of the open seek of the pervious for the open seek of the open seek of the pervious for the open seek of the open seek of the pervious for the open seek of the op

x86_64.tar.gr

Iding TopHat from source order to build TopHat2 you m

解析手順に関する記載がある 必要ツールなどの記載・ インストールの方法・ テストデータ等での極く簡単な

必要シール

- bowtie2
- samtools

ので、自分でmakeする必要はない。 バイナリーファイルが配布されている TopHat2はあらかじめコンパイルした 自分でソースからmakeする場合は

- SAMtools lib
- が必要 Boost C++ library

testデータが用意されている

tar zxvf test_data.tar.gz tophat -r 20 test_ref reads_1.fq reads_2.fq cd test_data

TopHat

read mapper for RNA-Sec



TopHat is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mand then analyzes the mapping results to identify splice junctions between exons. using the ultra high-thro

TopHat is a collaborative effort among Daehwan Kim and Steven Salzberg in the Center for Computational Biology at Johns Hopkins University, and Cole Trapnell in the Genome Sciences Department at the University of Washington. TopHat was originally developed by Cole Trapnell at the Center for Bioinformatics and Computational Biology at the University of Maryland, College Park.

TopHat 2.1.1 release 2/23/2016
TopHat 2.1.1 release 2/23/2016
Please note that TopHat has entered a low maintenance, low support stage as it is now largely superseded by HISAT2 which provides the same core functionality (i.e. spliced alignment of RNA-Seq reads), in a more accurate and much more efficient way.
Version 2.1.1 is a maintenance release which includes the following changes, some of them thanks to GitHub contributors:
TopHat can be now built on more recent Linux distributions with newer GNU C++ (5.x), as the included SeqAn library was finally upgraded to

- detection of linker options for the Boo st::Thread library which preve nted the TopHat build from source on so syst
- incorporated Luca Venturini's code to support large bowtie2 indexes (.bt2).

 bam2fastx usage message (-h/-help) was updated in order to better document the functions of this program which can be used as a standalone utility for converting reads from BAM/SAM to FASTQ/FASTA; the -v/-version option was also added to this utility for easier integration in other pipelines.

- TopHat 2.1.0 release 6/29/2015

 TopHat-Fusion algorithm improves ents for more sensitive and accurate discovery of fusions, thanks to contributions from Gordon Bean and
- pair-end insert) are counted as supporting evidence for the fusion parameter can be set using the --fusion-pair-dist <int> flag. fixed a few issues with GFF parsing of some annotation files fixed a runtime-error when using --no-discordant option. Ryan Kelley at Illumina.

 * This release implements a new algorithm for counting fusion-supporting read pairs that reduces the number of false-positive potential fusions. This algorithm computes the inner distance between read pairs by first converting the pair positions to transcript coordinates using the transcript information in refigene.txt and ensGene.txt. Pairs with small inner distance (suggesting the pair could come from a plausible pair-end insert) are counted as supporting evidence for the fusion. The default threshold for the inner distance is 250 base pairs; this

- Several fixes/improvements thanks to contributors on GitHub:

 new --max-num-fusions option allowing the user to specify the ma
 adjusting lower limit for --fusion-multipairs
 fixed a few bonn. number of reported fusions in tophat-fusion-post
- · fixed a few typos, cleaning up python code etc.

TopHat source code moved to GitHub 3/31/2015

bmit patches (pull requests). welcome to subm

- TopHat 2.0.14 release 3/24/2015 ersion 2.0.14 is a maintenance release
- ase with the following changes:
- pipeline speed improvements thanks to contributions from Véronique Legrand and M
 added support for xz compressed read files (thanks to a patch submitted by Ashton Trey Belew)

パラメータの意味など 詳しく知るためには、

必ずManualを見る ed to

Site Map

OSI certified

FAQ

News and updates

New releases and related tools will be announced through the Bowtie ailing list.

Getting Help

private communications only. Please do not email technical questions to TopHat contributors directly. TopHat can be posted on the **Tuxedo Tools Users Google Group.** Please use tophat.cufflinks@gmail.com for Questions and comments about

Releases

ion 2.1.1

2/23/2016

Related Tools

Linux x86_64 binary Mac OS X x86_64 binary

Cufflinks: Isoform assembly and quantitation for RNA-Seq Bowtie: Ultrafast short read alignment TopHat-Fusion: An algorithm for Discovery of Novel Fusion Transcripts CummeRbund: Visualization of RNA-

Manual

What is TopHat?

・リファレンス annotationなしでも ・75base以上のreadに最適化

スプライスジャンクションを見つける

- Using TopHat

What is TopHat?

TopHat is a program that aligns RNA-Seq reads to a genome in order to identify exon-exon splice junctions. It is built on the ultrafast short read mapping program Bowtie. TopHat runs on **Linux** and **OS X**.

What types of reads can I use TopHat with?

reads 75bp or longer. TopHat was designed to work with reads produced by the Illumina Genome Analyzer, although users have been successful in using TopHat with reads from other technologies. In TopHat 1.1.0, we began supporting Applied Biosystems' Colorspace format. The software is optimized for

How does TopHat find junctions?

possible splice junctions and then maps the reads against these junctions to confirm them. TopHat can find splice junctions without a reference annotation. By first mapping RNA-Seq reads to the genome, TopHat identifies potential exons, since many RNA-Seq reads will contiguously align to the genome. Using this initial mapping information, TopHat builds a database of

Short read sequencing machines can currently produce reads 100bp or longer but many exons are shorter than this so they would be missed in the initial mapping. TopHat solves this problem mainly by splitting all input reads into smaller segments which are then mapped independently. The segment alignments are put back together in a final step of the program to produce the end-to-end read alignments.

10 million). This latter option will only report alignments across "GT-AG" introns sequence or when an internal segment fails to map - again suggesting that such reads are spanning multiple exons. With this approach, "GT-AG", "GC-AG" and "AT-AC" introns will be found ab initio. The second source is pairings of "coverage islands", which are distinct regions of piled up reads in the initial mapping. Neighboring islands are often spliced together in the transcriptome, so TopHat looks for ways to join these with an intron. We only suggest users use this second option (--coverage-search) for short reads (< 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with a small number of reads (<= 45bp) and with TopHat generates its database of possible splice junctions from two sources of evidence. The first and strongest source of evidence for a splice junction is when two segments from the same read (for reads of at least 45bp) are mapped at a certain distance on the same genomic

inuminan has provided the RNA-Seq user community with a set of genome sequence indexes (including Bowtie indexes) as well as GTF transcript annotation files. These files can be used with TopHat and Cufflinks to quickly perform expression analysis and gene discovery. The annotation files are augmented with the rss_id and p_id GTF attributes that Cufflinks needs to perform differential splicing, CDS output, and promoter user analysis. We recommend that you download your Bowtie indexes and annotation files from this page. More information about Illumina's iGenomes project can be found here.

		Bos taurus			<u> </u>	c	Rattus norvegicus N	П	c	=	Mus musculus	2	Е	c	-	Tollo sabielly			Е	Organism D:
	NCBI			LISCHIO	nsembl	UCSC	NCBI	Ensembl	000	280	NCDI	Co	Ensembl	OCSC	Ce C		NCBI		Ensembl	Data source
UMD_3.1	Btau_4.6.1	Btau_4.2	Btau 4.2	UMD3.1	Btau_4.0	rn4	RGSC_v3.4	RGSC3.4	mm10	mm9	build37.2	build37.1	NCBIM37	hg19	hg18	build37.2	build37.1	build36.3	GRCh37	Version
13990 MB	13448 MB	13357 MB	13357 MB	14042 MB	13315 MB	13710 MB	14234 MB	13725 MB	14193 MB	14537 MB	15725 MB	15260 MB	14428 MB	21058 MB	17349 MB	21450 MB	15850 MB	15814 MB	17297 MB	Size
May 11 16:08	May 11 16:09	May 11 14:11	May 11 14:11	May 11 12:41	May 11 14:18	May 15 22:32	May 15 23:58	May 15 22:33	Jun 14 11:29	May 14 21:12	May 14 22:52	May 15 17:53	May 14 22:13	May 14 15:36	May 14 15:31	May 14 17:54	May 14 19:04	May 14 19:36	May 14 17:23	Last Modified

Site Map

News and updates

mailing list. New releases and related tools will be announced through the Bowtie

Getting Help

Questions and comments about TopHat can be posted on the Tuxedo Tools Users Google Group. Please use tophat.cuffinks@gmail.com for private communications only. Please TopHat contributors directly.

ersion 2.0.12 6/24/2014

Source code Linux x86_64 binary Mac OS X x86_64 binary

Related Tools

Cufflinks: Isoform assembly and quantitation for RNA-Seq Bowtie: Ultrafast short read alignn TopHat-Fusion: An algorithm for

いるので有効活用できる indexファイルやannotation メジャーな生物種では ファイル等が配布されて

NATURE PROTOCOLS | PROTOCOL



of RNA-seq experiments with TopHat and Cufflinks Differential gene and transcript expression analysis

Pimentel, Steven L Salzberg, John L Rinn & Lior Pachter Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold

Affiliations † Contributions † Corresponding author

Published online 01 March 2012 Nature Protocols 7, 562-578 (2012) | doi:10.1038/nprot.2012.016

★ Citation Reprints Rights & permissions Article metrics

Abstract • Accession codes • References • Author information

computer time for typical experiments and ~1 h of hands-on time transcriptome sequencing data and available computing resources but takes less than 1 d of quality visualizations of analysis results. The protocol's execution time depends on the volume of assembly, lists of differentially expressed and regulated genes and transcripts, and publicationand experts alike. The protocol begins with raw sequencing reads and produces a transcriptome skills, these tools assume little to no background with RNA-seq analysis and are meant for novices tool for visualizing RNA-seq analysis results. Although the procedure assumes basic informatics covers several accessory tools and utilities that aid in managing data, including CummeRbund, a protocol describes in detail how to use TopHat and Cufflinks to perform such analyses. It also ones, as well as compare gene and transcript expression under two or more conditions. This seq) data. Together, they allow biologists to identify new genes and new splice variants of known discovery and comprehensive expression analysis of high-throughput mRNA sequencing (RNAprincipled analysis software. TopHat and Cufflinks are free, open-source software tools for gene complexity of data from RNA-seq experiments necessitate scalable, fast and mathematically splice variants and quantify expression genome-wide in a single assay. The volume and Recent advances in high-throughput cDNA sequencing (RNA-seq) can reveal new genes and

protocol論文も出ている

ただし今となっては少し古い

Freeではない

tophat基本コマンド

and then analyzes the mapping results to identify splice junctions between exons. mammalian-sized genomes using the ultra high-throughput short read aligner **Bowtie**, TopHat is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to

> tophat -G gene.gtf -o out_dir genome read_1.fastq read_2.fastq

まずgtfに基づき、トランスクリプトにmapさせ、ゲノム位置として戻す。mapしないリードはゲノムから探す -G/--GTF <GTF/GFF3 file>

tophatの出力

prep_reads.info
align_summary.txt
deletions.bed
insertions.bed
junctions.bed
accepted_hits.bam
unmapped.bam

sam/bam フォーマットのファイル accepted_hits.bamファイルがこの後必要

账꾆1

tophatを用いてpreprocessing済みの2D_rep1.fastqファイルをgenome_chr4にmapさせよGTFファイルとしてgenes_chr4.gtfを用いる

例) これにファイルディレクトリーを加える tophat -p 4 -G genes_chr4.gtf -o 2D_rep1 genome_chr4 2D_rep1_R1.fastq 2D_rep1_R2.fastq 出力を確認しよう。

例えば、align_summary.txtを見ればどの程度mapしたか分かる。これでRNA-Segのリード配列がゲノム配列にアラインできた。

cufflinksを用いてアラインされたreadを数える

定義した方法でのカウントが可能 gene単位 トランスクリプト単位 エキソン単位

- cufflinks
- -BEDTools
- -HTseq
- が利用できる

今回はCufflinksを利用

そもそもTopHat > Cufflinksの解析系は同じ開発元、非常に良く使われている。

それに基づいて、genes単位、isoforms単位での解析を進めてくれる。 ローカスアノテーション情報を記載したgtfファイルを用意しておけば、

簡易的に、特定ローカスの解析などを進めたい場合や、gtfファイルがない場合などは、BEDToolsも有用gtfファイルを自分で作製するのは結構大変だが、bedファイルは比較的容易

http://cole-trapnell-lab.github.io/cufflinks/

Cufflinks

Transcriptome assembly and differential expression analysis for RNA-Seq.

Cufflinks assembles transcripts, estimates their abundances, and tests for differential expression and regulation in RNA-Seq samples. It accepts aligned RNA-Seq reads and assembles the alignments into a parsimonious set of transcripts. Cufflinks then estimates the relative abundances of these transcripts based on how many reads support each one, taking into account biases in library preparation protocols.

Cufflinks was originally developed as part of a collaborative effort between the Laboratory for Mathematical and Computational Biology, led by Lior Pachter at UC Berkeley, Steven Salzberg's computational genomics group at the Institute of Genetic Medicine at Johns Hopkins University, and Barbara Wold's lab at Caltech. The project is now maintained by Cole Trapnell's lab at the University of Washington.

Cufflinks is provided under the OSI-approved Boost License

News

To get the latest updates on the Cufflinks project and the rest of the "Tuxedo tools", please subscribe to our mailing list

MARCH 25, 2014 APRIL 11, 2013	Cufflinks 2.2.0 released Cufflinks 2.1.1 released
MAY 05, 2014	Cufflinks 2.2.1 released
DECEMBER 10, 2014	Cufflinks has moved to GitHub

Protocol

switching during cell differentiation reveals unannotated transcripts and isoform Transcript assembly and quantification by RNA-Seq

Cole Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J van Baren, Steven L Salzberg, Barbara J Wold & Lior Pachter

Affiliations | Contributions | Corresponding author

Received 02 February 2010 | Accepted 22 March 2010 | Published online 02 May 2010 Nature Biotechnology 28, 511-515 (2010) | doi:10.1038/nbt.1621



discovery and abundance estimation 1,2,3 . However, this would require algorithms that are based genome annotation. even this well-studied model of muscle development and that it can improve transcriptomesuggest that Cufflinks can illuminate the substantial regulatory flexibility and complexity in or splice isoform, and we observed more subtle shifts in 1,304 other genes. These results series, 330 genes showed complete switches in the dominant transcription start site (TSS) independent expression data or by homologous genes in other species. Over the time known transcripts and 3,724 previously unannotated ones, 62% of which are supported by reads from a mouse myoblast cell line over a differentiation time series. We detected 13,692 Cufflinks. To test Cufflinks, we sequenced and analyzed >430 million paired 75-bp RNA-Seq splicing. Here we introduce such algorithms in an open-source software program called not restricted by prior gene annotations and that account for alternative transcription and High-throughput mRNA sequencing (RNA-Seq) promises simultaneous transcript



Cufflinks

Transcriptome assembly and differential expression analysis for RNA-Seq.

Cufflinks is available for Linux and Mac OS X. You can find the full list of releases below.

grab the current code, check out the Cufflinks GitHub repository. The Cufflinks source code for each point release is available below as well. If you want to



Cufflinks Releases

Date

2.1.0	2.1.1	2.2.0	2.2.1
April 10, 2013	April 11, 2013	March 25, 2014	May 05, 2014
Linux	Linux	Linux	Linux
Mac OS X	Mac OS X	Mac OS X	Mac OS X
Source	Source	Source	Source

Cufflinks

Transcriptome assembly and differential expression analysis for RNA-Seq.

- Install quick-start
- Installing a pre-compiled binary release
- **Building Cufflinks from source**
- Installing Boost
- Installing the SAM tools
- Installing the Eigen libraries

が必要

Boost C++ library

自分でソー Samtools

スからmakeする場合は

- Building Cufflinks
- Testing the installation
- Common uses of the Cufflinks package
- Using pre-built annotation packages

cufflinks ./test_data.sam

これでツールが動くことを確認

Install quick-start

Installing a pre-compiled binary release

a directory in your PATH environment variable. your machine, untar it, and make sure the cufflinks,cuffdiff and cuffcompare binaries are in from occasionally frustrating process of building Cufflinks, which requires that you install the Boost libraries. To use the binary packages, simply download the appropriate one for In order to make it easy to install Cufflinks, we provide a few binary packages to save users

INSTALL MANAUL GETTINGSTARTED TOOLS WELD HOWITHOOMS PROTOCOL BENCHMARKS CODE BAFED

Cufflinks

ranscriptome assembly and differential expression analysis for RNA-Seq.

Bowtie: ultrafast short read alignment

Bowlie is an ultrafast and memory-efficient tool for aligning sequencing reads to long reference sequences. It is particularly good at aligning reads of about 50 up to 108s or 1,000s of characters, and particularly good at aligning to relatively long (e.g., mammalian) genomes. Bowlie 2 indexes the genome with an FM index to keep its memory footprint small: for the human genome, its memory footprint is typically around 3.2 GB. Bowtie 2 supports gapped, local, and paired-end alignment modes.

Bowtie is provided under the OSI-approved Artistic License 2.0.

TopHat: alignment of short RNA-Seq reads

TopHat is a fast splice junction mapper for RNA Seq reads. It aligns RNA Seq reads to mammalian sized genomes using the ultra high-throughput short read aligner Bowtie, and then analyzes the mapping results to identify splice junctions between exons.

TopHat is provided under the OSI-approved Artistic License 2.0.

CummeRbund: visualization of RNA-Seq differential analysis

CummeRbund is an R package that is designed to aid and simplify the task of analyzing Cufflinks RNA-Seq output.

CummeRbund is provided under the OSI-approved Artistic License 2.0.

Monocle: Differential expression for singlecell RNA-Seq and qPCR.

Monocle is a toolkit for analyzing single-cell gene expression experiments. Monocle was designed for RNA-Seq, but can also work with single cell qPCR. It performs differential expression analysis, and can find genes that differ between cell types or between cell states. When used to study an onganig biological process such as cell differentiation, Monocle learns that process and places cells in order according to their progress through it Monocle finds genes that are dynamically regulated during that process.

Monocle is provided under the OSI-approved Artistic License (version 2.0)

Cufflinksの関連ツール Bowtie, TopHatは説明済み

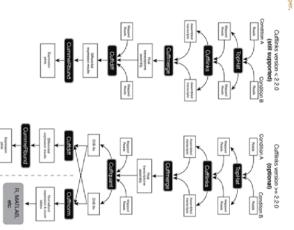
INSTALL MANUAL GETTING STAFTED TOOLS HELP HOWIT WOODS PROTOCOL BENCHMADES CODE NATE

Cufflinks

anscriptome assembly and differential expression analysis for RNA-Seq

The Cufflinks RNA-Seq workflow

The Cufflinks suite of tools can be used to perform a number of different types of analyses for RRA-Seq experiments. The Cufflinks suite includes a number of different programs that work together to perform these analyses. The complete workflow, performing all the types of analyses Cufflinks can execute, is summarized in the graph below. The left side illustrate the "classie" RRA-Seq workflow, which includes read mapping with Top-Nat, assembly with Cufflinks, and visualization and exploration of results with CummeBound. A newer, more advanced workflow was introduce with Cufflinks version 2.2.0, and is shown on the right. Both are still supported. You can read about the classic workflow in detail in our protocol



Cufflinks

Cufflinks is both the name of a suite of tools and a program within that suite. Cufflinks the program assembles transcriptomes from RNA-Seq data and quantifies their expression.

Cuffcompare

After assembling a transcriptome from one or more samples, you'll probably want to compare your assembly to known transcripts. Even if there is no "reference" transcriptome for the organism you're studying, you may want to compare the transcriptomes assembled from different RNA-Seq libraries, culfcompare helps you perform these comparisons and assess the quality of your assembly.

Cuffmerge

When you have multiple RNA-Seq libraries and you've assembled transcriptomes from each of them, we recommend that you merge these assemblies into a master transcriptome. This step is required for a differential expression analysis of the new transcripts you've assembled. Cuffmerge performs this merge step.

Cuffquant

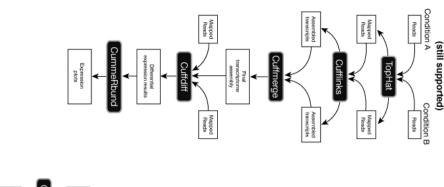
Quantifying gene and transcript expression in RNA-Seq samples can be computationally expensive. Culfquant allows you to compute the gene and transcript expression profiles and save these profiles to files that you can analyze later with Cuffdiff or Cuffnorm. This can help you distribute your computational load over a cluster and is recommended for analyses involving more than a handful of libraries.

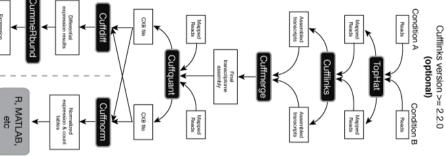
Cuffdiff

Comparing expression levels of genes and transcripts in RNA-Seq experiments is a hard problem. Cuffdiff is a highly accurate tool for performing these comparisons, and can tell you not only which genes are up- or down-regulated between two or more conditions, but also which genes are differentially spliced or are undergoing other types of isoform-level regulation.

Cuffnorm

Sometimes, all you want to do is normalize the expression levels from a set of RNA-Seq libraries so that they're all on the same scale, facilitating downstream analyses such as clustering. Expression levels reported by Cufflinks in FPKA units are usually comparable between samples, but in certain situations, applying an extra level of normalization can remove sources of bias in the data. Cuffnorm normalizes a set of samples to be on as similar scales as possible, which can improve the results you obtain with other downstream





cufflink cufflinks

Cufflinks version < 2.2.0

cuffmerge cuffcompare

cuffquant cuffnorm cuffdiff

の6つのプログラムから構成

cuffquant, cuffnormは ver2.2.0(20140325) から実装

MacOSX版のバイナリーはver2.2.0以降はバグがありsegmentation errorでまともに動かないようです。

今回の実習ではver2.1.1を使用し、 cuffquant, cuffnormは簡単な説明のみ に留めます。

Cufflinks

analysis for RNA-Seq

and isoform switching during cell RNA-Seq reveals unannotated transcripts Transcript assembly and quantification by differentiation

どうやって動いているか

なったら詳緒を拍描していく まず動いて使えそうな感じに

Nature Biotechnology, 2010

y prior gene annotations and that account for alternative transcript discovery trior gene annotations and that account for alternative transcription and splicing. Here the introduce such algorithms in an open-source software program called Cufflinks. To test ufflinks, we sequenced and analyzed >430 million paired 75-bp RNA-Seq reads from a ouse myoblast cell line over a differentiation time series. We detected 13,692 known anscripts and 3,724 previously unannotated ones, 62% of which are supported by dependent expression data or by homologous genes in other species. Over the time ries, 330 genes showed complete switches in the dominant transcription above a splice isoform, and we observed more witchigh-throughput mRNA sec erved more subtle shifts in 1,304 other genes. These results minate the substantial regulatory flexibility and complexity in of muscle development and that it can improve annotation.

i:10.1038/nbt.1621

by correcting for fragment bias Improving RNA-Seq expression estimates

Genome Biology, 2011

cufflinks基本コマンド

Cufflinksコマンド

cufflinks -o out_directory -G hoge.gtf tophat_directory/accepted_hits.bam

cufflinksを実行してパラメー -タを確認しよう。

考慮す べきパラメー -女一例

- 出力の指定、TopHatの出力と同じ場所にしておくのが分かりやすいだろうCPUスレッド数の指定(デフォルトは1)、結構時間がかかるので使える数を指定すると良いだろうGTFファイルに記載されたアノテーションのみについて計算GTFファイルに記載されたアノテーションをガイドにしてアセンブルする無視したいトランスクリプト(rRNAなど)を指定
- Фф 9
- \$ מַּ

cufflinks出力

skipped.gtf
transcripts.gtf
genes.fpkm_tracking
isoforms.fpkm_trancking

実習2 先のtophatの結果を用いてcufflinksにかけてみよう

出力を確認しよう。 > cufflinks -p 4 -o 例) これにファイルディフクトリーを加える 2D_rep1 -G genes_chr4.gtf accepted_hits.bam

また-gを用いてcufflinksにかけると新規の発現領域が存在するのが分かる geneごと、isoformごとにFPKM値が計算されているのが分かる。

cuffcompareコマンド

you assemble. The program cuffcompare helps you: Cufflinks includes a program that you can use to help analyze the transfrags

Compare your assembled transcripts to a reference annotation Track Cufflinks transcripts across multiple experiments (e.g. across a time course)

From the command line, run cuffcompare as follows:

cuffcompare [options]* <cuff1.gtf> [cuff2.gtf] ... [cuffN.gtf]

今回はすでにあるgtfファイルの情報を用いるので、意識的に使う必要はない。

cuffmergeコマンドと出力

個々のサンプルのアセンブルモデルを統合する。

```
Usage: cuffmerge [Options] <assembly_GTF_list.txt>
                                                           Options:
-g/--ref-gtf
-s/--ref-sequence
--min-isoform-fraction
-p/--num-threads
--keep-tmp
        [ default:
           . 05
1
```

統合ファイルリストを事前に作製する必要がある(例 assemblies,txt)

cuffmerge -s \$REFSEQ -g \$GTF assemblies.txt

例 assemblies,txt

~/arabi_2D_2/transcripts.gtf ~/arabi_2D_3/transcripts.gtf ~/arabi_2D2L_2/transcripts.gtf ~/arabi_2D2L_3/transcripts.gtf

Cufflinks includes a script called cuffmerge that you can use to merge together several Cufflinks assemblies. It handles also handles running Cuffcompare for you, and automatically filters a number of transfrags that are probably artifiacts. If you have a reference GTF file available, you can provide it to the script in order to gracefully merge novel isoforms and known isoforms and maximize overall assembly quality. The main purpose of this script is to make it easier to make an assembly GTF file suitable for use with Cuffdiff.

± ±

merged.gtf

今回はすでにあるgtfファイルの情報を用いるので、 使う必要はない。

cuffdiffコマンド

DE gene等を統計計算で取り出す コマンド入力して使用法を確認してみよう

```
Usage: cuffdiff [options] <transcripts.gtf> <sample1_hits.sam> <sample2_hits.sam> [... sampleN_hits.sam]
Supply replicate SAMs as comma separated lists for each condition:
sample1_rep1.sam,sample1_rep2.sam,...sample1_repM.sam
General Options:
-o/--output-dir
-t/--labels

Comma-separated i:...sample1_hits.sam> [... sample2_hits.sam> [... sampleN_hits.sam]
comma-separated list of condition la
False discovery rate used in testing
   default:
   0.05 ]
```

cuffdiff -o out_file merged.gtf bam1,bam2,bam3 bam4,bam5,bam6

cuffdiffにかかる時間やメモリー使用量が軽減される。 Version 2.2.0以降は後述のcuffquantで得られたcxbファイルをbamファイルの代わりに用いる。

cuffdiffの出力

bias_params.info
run.info
read_groups.info
var_model.info
cds.read_group_tracking
cds.fpkm_tracking
cds.count_tracking
genes.read_group_tracking
genes.count_tracking
isoforms.read_group_tracking
isoforms.read_group_tracking
isoforms.fpkm_tracking
isoforms.fpkm_tracking
tss_groups.read_group_tracking
tss_groups.read_group_tracking

gene_exp.diff
cds_exp.diff
cds.diff
isoform_exp.diff
promoters.diff
splicing.diff
tss_group_exp.diff

diffの付いたファイルがそれぞれの 違いの情報を記載したファイル

.diffファイルの内容

13	12	11	10	9	80	7	6	5	4	ω	2	-	Column
significant	q value	p value	test stat	log2 (FPKM _y /FPKM _x)	FPKM _y	FPKM _x	Test status	sample 2	sample 1	locus	gene	Tested id	Column name
no	0.985216	0.389292	0.860902	0.06531	8.551545	8.01089	NOTEST	Brain	Liver	chr1:4797771- 4835363	Lypla1	XLOC_000001	Example
Can be either "yes" or "no", depending on whether p is greater then the FDR after Benjamini-Hochberg correction for multiple-testing	The FDR-adjusted $ ho$ -value of the test statistic	The ${\it uncorrected}\ p$ -value of the test statistic	The value of the test statistic used to compute significance of the observed change in FPKM	The (base 2) log of the fold change y/x	FPKM of the gene in sample y	FPKM of the gene in sample x	Can be one of OK (test successful), NOTEST (not enough alignments for testing), LOWDATA (too complex or shallowly sequenced), HIDATA (too many fragments in locus), or FAIL, when an ill-conditioned covariance matrix or other numerical exception prevents testing.	Label (or number if no labels provided) of the second sample being tested	Label (or number if no labels provided) of the first sample being tested	Genomic coordinates for easy browsing to the genes or transcripts being tested.	The gene_name(s) or gene_id(s) being tested	A unique identifier describing the transcipt, gene, primary transcript, or CDS being tested	Description

cuffquantコマンドと出力(ver2.2.0以降)

bamの内容からgene/transcriptレベルで定量化し、バイナリー出力する

cuffquant -o out_directory hoge.gtf accepted_hits.bam

cuffquantを実行してパラメータを確認しよう。

考慮すべきパラメ-

- -o 出力ディレクトリーの指定 -p CPUスレッド数の指定(デフォルトは1)、結構時間がかかるので使える数を指定 -M 無視したいトランスクリプト(rRNAなど)を指定 他にもestimationに関わる -b -uパラメータがある。

出力

abundances.cxb

> cuffquant -p 4 -o 2D_1 genes_chr4.gtf accepted_hits.bam

出力ファイルはこの1つだけ 新たにcxbファイルが作製されていることが分かる。

cuffdiffを速くできる。 cuffdiffの前にcuffquantを行い、cxbファイルを作製することで

cuffnormコマンドと出力(ver2.2.0以降)

Cuffnormコマンド

Cuffnorm, which simply computes

a normalized table of expression values for genes and transcripts.

> cuffnorm -o out_file genes_chr4.gtf bam1,bam2,bam3 bam4,bam5,bam6

[sampleN.sam_replicate1.sam[,...,sample2_replicateM.sam]] cuffnorm [options]* <transcripts.gtf> <sample2_replicate1.sam[,...,sample2_replicateM.sam]>... <sample1_replicate1.sam[,...,sample1_replicateM.sam]>

sam/bamかcxbファイルどちらも入力可能。 ただし混在は不可

cuffnormの出力(ver2.2.0以降)

genes.count_table genes.attr_table cds.count_table tss_groups.fpkm_table tss_groups.count_table samples.table run.info isoforms.count_table isoforms.attr_table genes.fpkm_table cuffnorm.tree cds.fpkm_table cds.attr_table isoforms.fpkm_table _groups.attr_table

たくさんのサンプルで発現プロットやクラスター図を書きたい場合便利。

tophat -> cufflinksの解析系を使用する際の注意

or plot expression levels of genes important in samples and you simply want to cluster them output files are useful when you have many analysis. To assess the significance of changes your study. between conditions, use Cuffdiff. Cuffnorm's in expression for genes and transcripts It does not perform differential expression

differential expression tools that require raw size, they should not be used with downstream that because these counts are already transcript, TSS group, and CDS group. Note fragments that originate from each gene, Cuffnorm will report both FPKM values and counts as input. normalized to account for differences in library normalized, estimates for the number of

tophat -> cufflinksは一連の解析系

cufflinksの出力はすでにノ inputには利用できない。 するedgeRなどの別のツー されたものなので、rawデ 670 ・マリイズ

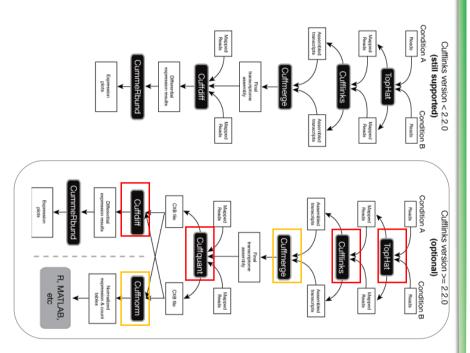


Count-based differential expression analysis of RNA sequencing data using R and Bioconductor

Mocols 8, 1765–1786 (2013) | doi:10.103 online 22 August 2013

full text @ PDF & Citation lin Rep M An

versionによる違いまとめ



tophat, cufflinksの実習

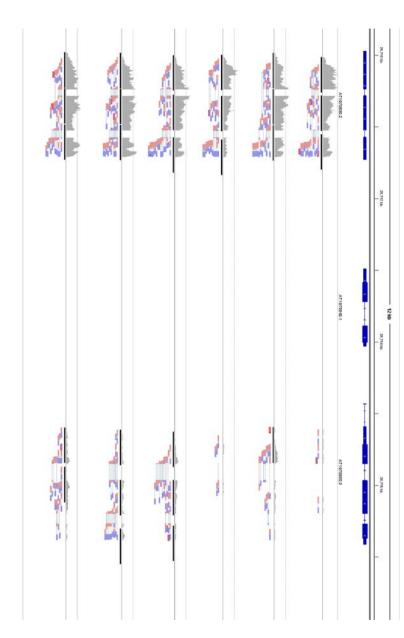
/data/KY/tophatにあるファイルを用いる。

1. TopHatを用いて、paired-endのtest data(2D_rep1_R1.fastq,2D_rep1_R2.fastqをリファレンスgenome_chr4にマップさせよオプション -Gの有無に よる違いを確認しよう。

2.Cufflinksを用いて、 2D_rep1のカウントをしよう。 -Gと-gの違いを確認しよう。

RNA-Seq結果のIGV実習

TAIR10の配列を呼び出し、TopHatで得られたBAMファイルをindexファイルを付け、読み込む map結果をIGVで可視化してみよう



Excelを使って結果を確認してみよう

Excelでgene_exp.diffファイルを読み込んでみるtab区切りテキストファイルなのでそのまま読み込める Excelのsort機能を使ってq値でsortしてみる 2D vs 2D2Lのcuffdiff結果が~data/KY/cuffdiffフォルダーにある。



test_id	gene_id	gene	locus	sample_1	sample_2	status	value_1	value_2	log2(fold_chan test_stat p_value	test_stat		q_value
XLOC_000047	XLOC_000047	KEA1	1:284609-291094	q1	q2	OK.	12.8356	47.6879	1.89347	4.44122	5.00E-05	0.000325 yes
XLOC_000091	XLOC_000091	BXL2	1:564204-567769	q1	q2	OK	112.839	21.5634	-2.38762	-6.02938	5.00E-05	0.000325 yes
XLOC_000148	XLOC_000148	PSB27	1:898875-899655	q1	q2	OK	194.744	691.64	1.82844	7.10401	5.00E-05	0.000325 yes
XLOC_000310	XLOC_000310	PSBP-1	1:2047824-2049418	q1	q2	OK	588.195	3147.84	2.42	7.92975	5.00E-05	0.000325 yes
XLOC_000404	XLOC_000404	NPQ1	1:2706923-2709531	q1	q2	OK.	21.2494	78.5734	1.88662	3.26377	5.00E-05	0.000325 yes
XLOC_000419	XLOC_000419	CSD1	1:2827060-2838469	q1	q2	OK.	503.523	181.545	-1.47173	-5.38312	5.00E-05	0.000325 yes
XLOC_000450	XLOC_000450	CSP41B	1:3015327-3018234	q1	q2	OK	113.687	650.406	2.51627	8.83387	5.00E-05	0.000325 yes
XLOC_000487	XLOC_000487	LRR XI-23	1:3252239-3255693	q1	q2	Q.	26.4081	49.6396	0.910512	2.30664	5.00E-05	0.000325 yes
XLOC_000598	XLOC_000598	ATGLX1	1:3995168-3997907	q1	q2	Q	60.1583	162.387	1.4326	3.26419	5.00E-05	0.000325 yes
XLOC_000600	XLOC_000600	AT1G11860	1:4001112-4003442	q1	q2	OK	319.6	756.582	1.24323	4.18318	5.00E-05	0.000325 yes
XLOC_000614	XLOC_000614	AT1G12080	1:4084161-4085045	q1	q2	OK.	1884.29	67.9613	-4.79316	-9.20293	5.00E-05	0.000325 yes
XLOC_000616	XLOC_000616	CHL1-1	1:4105232-4109545	q1	q2	OK.	107.267	57.7917	-0.892267	-2.70294	5.00E-05	0.000325 yes
XLOC_000624	XLOC_000624	AT1G12230	1:4147961-4151056	q1	q2	Q	102.049	50.9296	-1.00268	-2.40566	5.00E-05	0.000325 yes
XLOC_000680	XLOC_000680	CYP71B7	1:4467219-4469033	q1	q2	Q	17.1443	84.588	2.30272	4.53043	5.00E-05	0.000325 yes
XLOC_000724	XLOC_000724	AT1G13930	1:4761011-4762666	Q1	q2	Q.	94.6747	2483.48	4.71324	10.4968	5.00E-05	0.000325 yes
XLOC_000749	XLOC_000749	AT1G14345	1:4899144-4899979	9.	q2	OK.	38.3992	157.145	2.03295	4.49341	5.00E-05	0.000325 yes
XLOC_000765	XLOC_000765	AT1G14670	1:5037611-5040528	q1	q2	OK.	84.8105	44.439	-0.932415	-2.66978	5.00E-05	0.000325 yes
XLOC_000835	XLOC_000835	NDF1	1:5489297-5493772	ę,	q2	e	20.0548	104.567	2.3824	4.27443	5.00E-05	0.000325 yes
XLOC_000884	XLOC_000884	HCF173	1:5723087-5727312	q1	q2	Q.	7.34039	112.227	3.93442	5.2414	5.00E-05	0.000325 yes
XLOC_000916	XLOC_000916	FUG1	1:5885082-5890470	Q1	q2	Q	48.9638	105.457	1.10687	3.5512	5.00E-05	0.000325 yes
XLOC_001003	XLOC_001003	NDF6	1:6460597-6462224	9.	q2	OK.	45.3045	185.555	2.03412	2.97075	5.00E-05	0.000325 yes
XLOC_001030	XLOC_001030	LHCA6	1:6612748-6613972	q1	q2	Q.	52.6816	153.395	1.54188	4.09397	5.00E-05	0.000325 yes
XLOC_001063	XLOC_001063	PUP14	1:6832346-6833837	q1	q2	Q	37.731	91.5568	1.27892	3.13218	5.00E-05	0.000325 yes
XLOC_001076	XLOC_001076	ATLFNR2	1:6942716-6945018	q1	q2	Q.	87.7487	1025.37	3.54662	10.0816	5.00E-05	0.000325 yes
XLOC_001099	XLOC_001099	AT1G20390	1:7065493-7071561	9.	q2	OK.	45.6232	15.9769	-1.51378	-4.22277	5.00E-05	0.000325 yes
XLOC 001170	XLOC 001170	AT1G21680	1:7613004-7615339	91	g2	읒	27.146	80 96	1 57647	2 02821	5 00F-05	0.000335 vas

GTFファイルに記載された遺伝子ごとの発現カウントに対して倍率、p値、q値が計算される。

Rを使ってMA plotを書いて見よう

先と同じgene_exp.diffファイルを読み込んでみるtab区切りテキストファイルなのでread.delim関数で読み込む colorのパラメplot関数を使って描画 M, Aをそれぞれ計算する -タをsignifitureの値で色分けさせてみる。

例)

plot(A,M,col=dat\$significant, pch=16, cex=0.4, ylim=c(-8,8)) M<-log2(dat\$value_1+1)-log2(dat\$value_2+1) $A<-1/2*(log2(dat$value_1+1)+log2(dat$value_2+1))$ dat <- read.delim("gene_exp.diff")

簡易スクリプトを使って、結果を成形してみよう

またその数を数えよ。 Awkによる1行スクリプトで、q_value値が0.05以下となる行を取り出せ。

廻)

q_valueが0.05以下のもののみリストアップするには? q_valueの記載は13列目だから・・・

awk '\$13<=0.05 {brint \$0}, gene_exp.diffと記述すればOK \$で列番号を指定できる \$0は行全体を意味する

awk '\$13<=0.05 {print \$0}' gene_exp.diff | wc で数も分かる。

実践演習課題

DE gene等を調べよ。 2D_rep*_R*.fastq(2days dark条件で育てたアラビドプシス芽生え), 2D2L_rep*_R*.fastq(その後2days light条件で育てたアラビドプシス芽生え) それぞれ3反復のデータ を用い TopHat→Cufflinksの系を用い、 cutadapt済みのデータセット~data/KY/tophatフォルダーの

GTFファイルとしてgenes_chr4.gtf fastaファイルとしてgenome_chr4.fa を利用する。 (アラビドプシスTAIR10の配列だが計算時間を考慮して、それぞれChr4のみになっている)

RNA-Seqパイプライン -ゲノムベースの解析法-の最後3スライドを参考に、マッピングデータのIGVでの可視化、

エクセルでの確認、

Rを用いたM-A plotの描画、

例に挙げた流れで、簡易スクリプトを用いたデ--夕描出を中よ。