

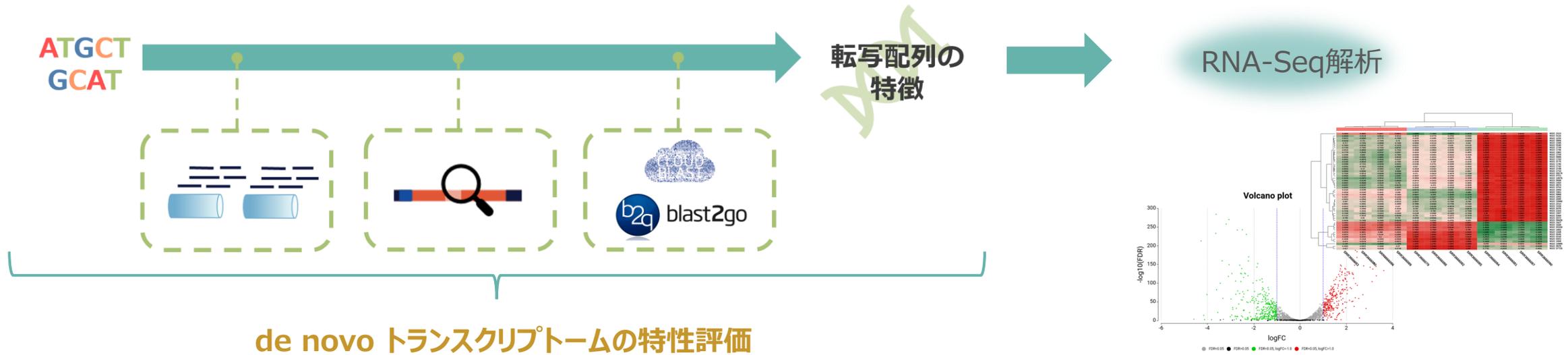


de novo トランスクリプトームの特性評価

フィルジェン株式会社

バイオインフォマティクス部(biosupport@filgen.jp)

- de novo トランスクリプトームはリファレンスゲノムの無い生物種やアノテーションの無い新規転写産物の探索に有用な解析アプローチです。
- この解析により、タンパク質の機能推定を得るだけでなく、リファレンスゲノムの無いあるいは不十分な生物種に対して、発現値定量や発現比較などのRNA-Seq解析を行うことができます。



QC・トリミング

- NGSより出力された生データが良好か、下流分析に影響する問題がないか確認。
RNA-Seq de novo アセンブリに有用な高品質なリードデータが得られる。
FastQCとTrimmomaticツールを統合

RNA-Seq de novo アセンブリ

- 高品質なリードデータのみで（リファレンスゲノムなしで）新規にゲノム配列を構築する解析。
構築された長い連続的な配列（コンティグ）に関するFASTAファイルが得られる。
Trinityツールを統合

コンティグの完全性評価

- RNA-Seq de novo アセンブリで作成したトランスクリプトーム配列の品質を確認する。
BUSCOを統合

配列クラスタリング

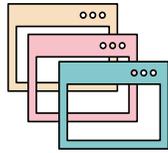
- 配列類似性に基づいてクラスタリングを行うことでトランスクリプトーム配列冗長性を取り除く。
CD-HITを統合

コーディング領域予測

- トランスクリプトーム配列からコード領域を探す。
TransDecoderを統合

アノテーション

- BlastやInterProScanにより遺伝子機能情報を付与することができる



- コマンドライン型のツールかつパラメータ設定が重要なため操作が煩雑



- RNA-Seq de novo アセンブリやBlast解析は、高スペックPCやサーバーの導入が必要



OmicsBoxのde novo トランスクリプトームの特性評価

- メーカーのサーバーで高速計算 高価なPCの購入は不要
- マウス操作&既成のワークフローで簡単に解析
- de novo トランスクリプトームの特性評価を行える数少ないソフトウェア

必要なファイル

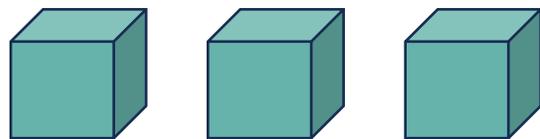


シーケンサーから出力された生データ
もしくは受託サービスで得られたクリーンリードデータ

例えば以下のような実験を行うとします：

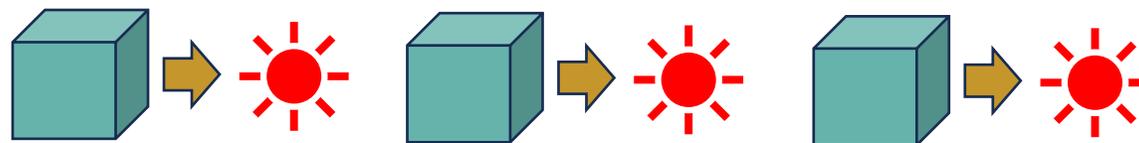
果実の褐色腐敗病の原因となる *Monilinia laxa*

暗所で 4 日間生育した菌糸体



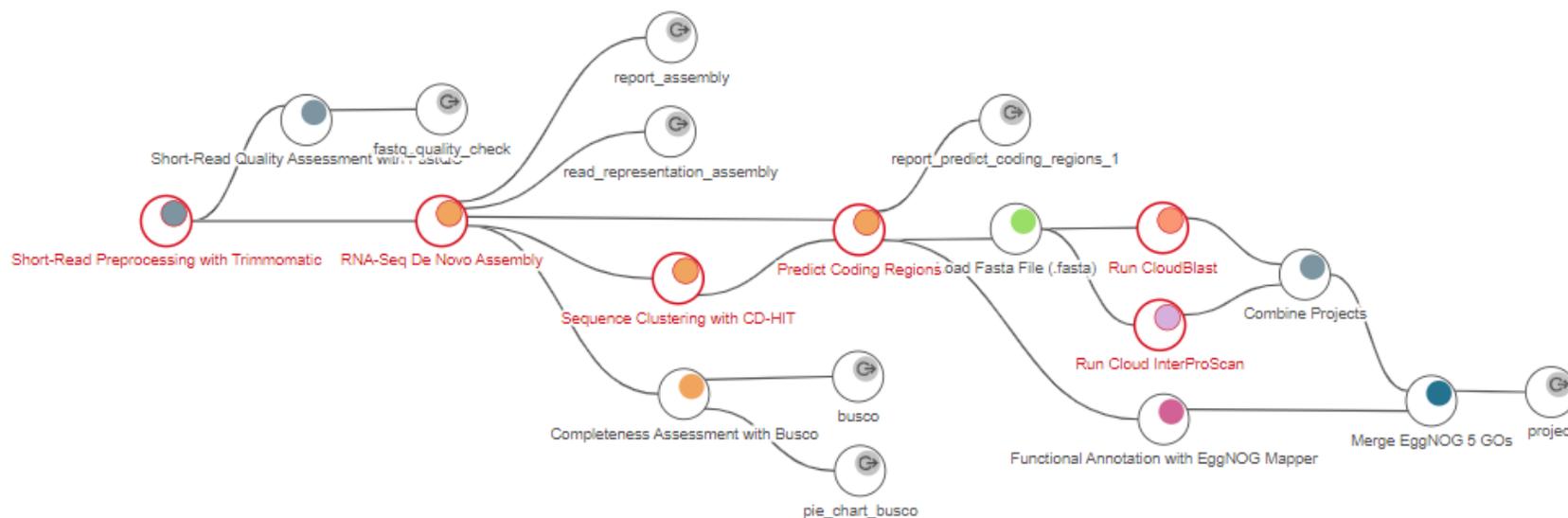
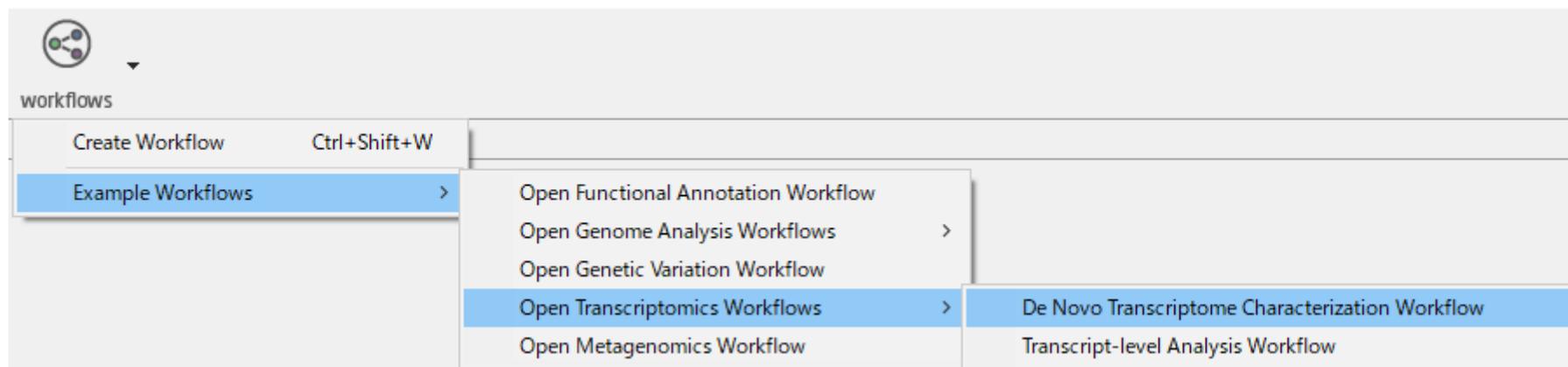
n=3

暗所で 2 日間生育し、2 日間光に曝露した菌糸体



n=3

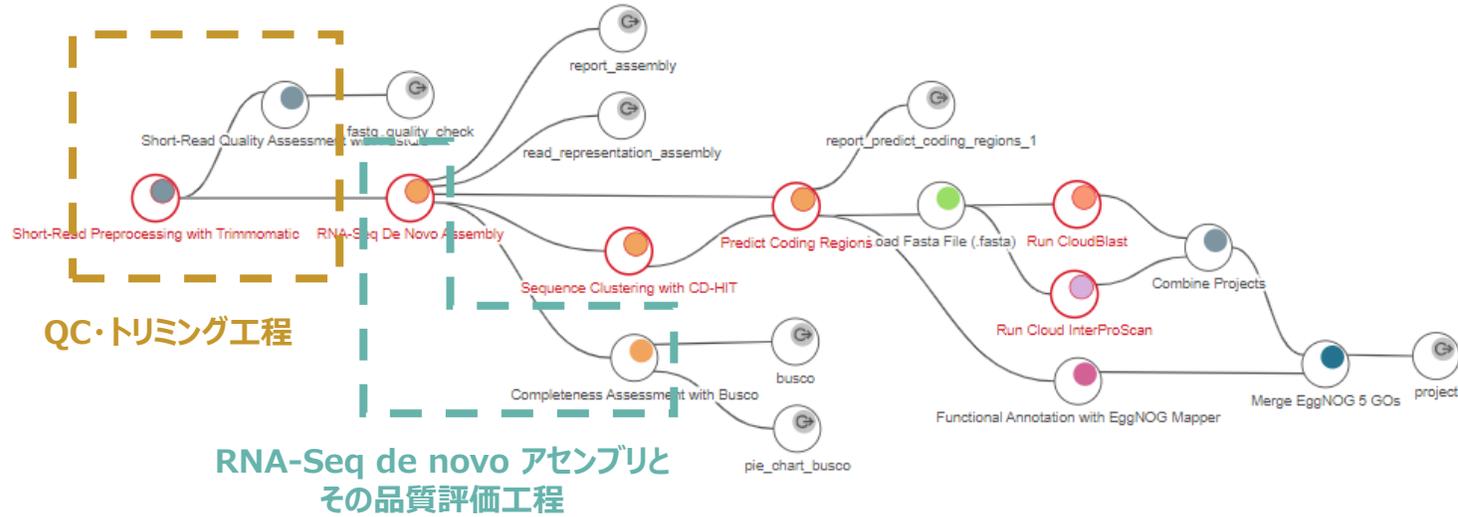
この場合、すべてのサンプル（計6サンプル）をまとめて解析し
1つのトランスクリプトーム配列を構築することになります。



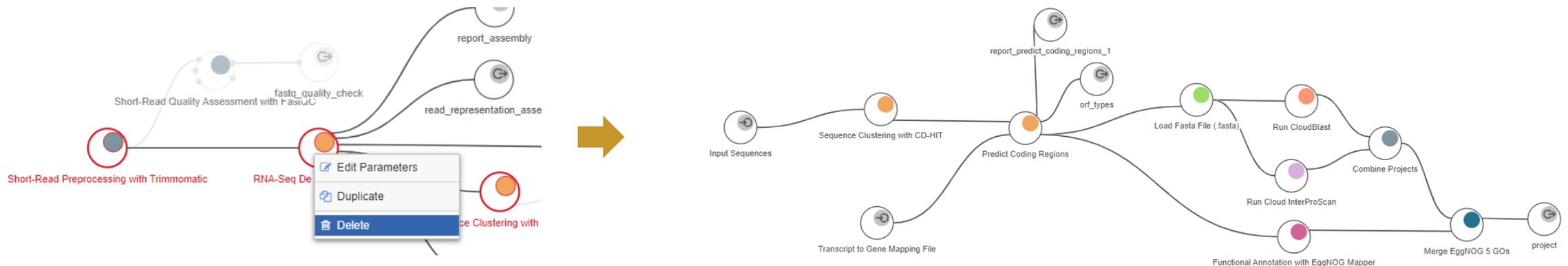
この解析のための既成のワークフローが搭載されています。

*Transcriptomics ModuleとFunctional Analysis Moduleを使用します。

OmicBoxのワークフロー機能

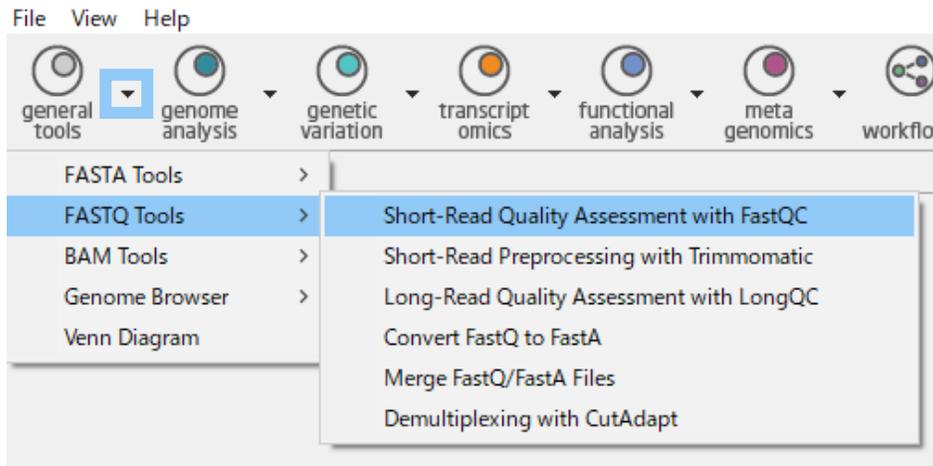


ただし、QC・トリミングとRNA-Seq de novo アセンブリとその品質評価については結果によっては、この解析を繰り返す必要があるため、これらの工程は個別に解析を行い以下のようにワークフローを修正できます。

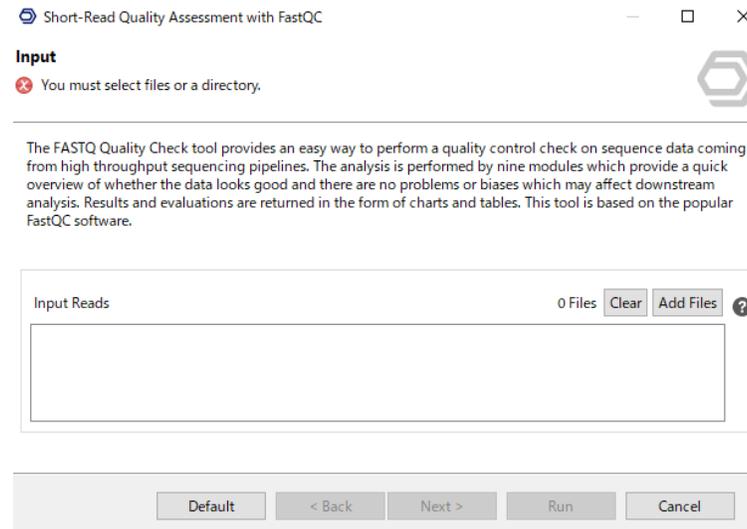




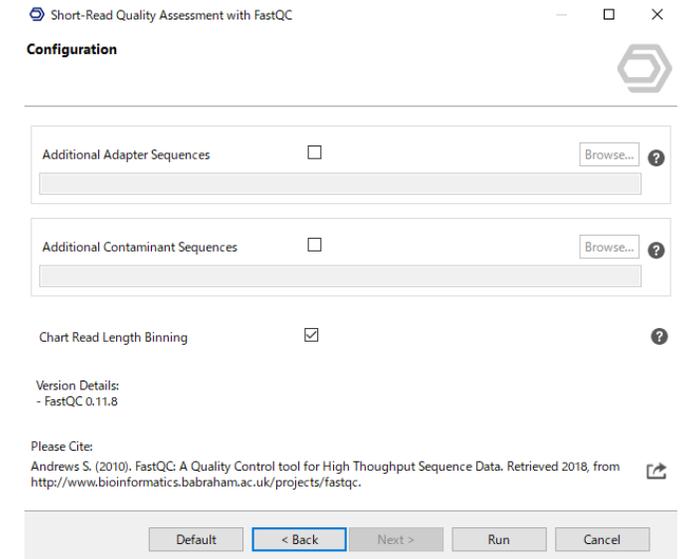
- データが良好か、下流分析に影響する問題がないか確認



画面上部のアイコンからQCツールを起動



全てのサンプルをこちらに指定



任意でアダプタリストを追加

結果

Welcome Message | FASTQ Quality Check (Dataset) | FASTQ Quality Check (ERR1948631_1.fastq) | FASTQ Quality Check (clean_ERR1948631_1.fq) | Chart: Adapter Content

FASTQ Quality Check

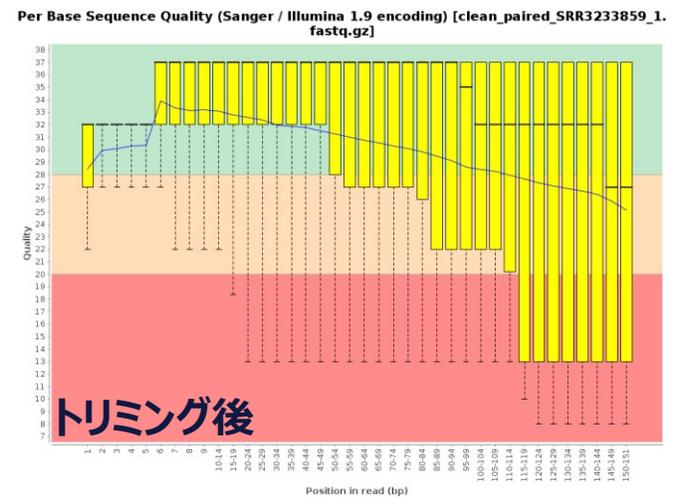
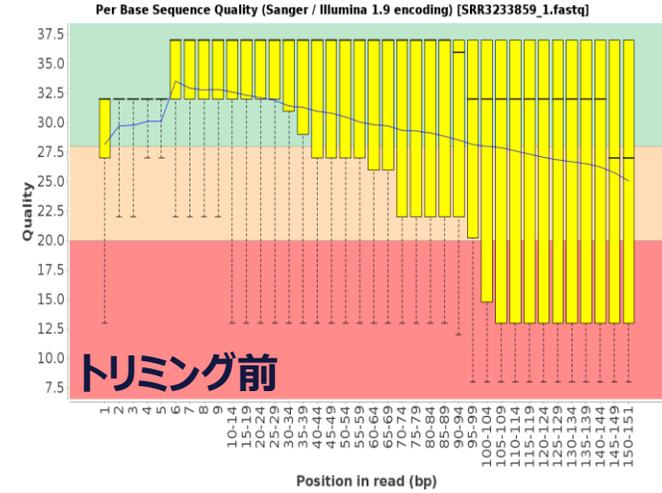
Name: Dataset

Overall Results

| Name | Per Base Sequence Quality | Per Sequence Quality Scores | Per Base Sequence Content | Per Sequence GC Content | Per Base N Content |
|-----------------------|---------------------------|-----------------------------|---------------------------|-------------------------|--------------------|
| ERR1948631_1.fastq | PASS | PASS | FAIL | PASS | PASS |
| clean_ERR1948631_1.fq | PASS | PASS | FAIL | PASS | PASS |

| Name | Sequence Length Distribution | Adapter Content | Overrepresented Sequences | Sequence Duplication Levels | Report |
|-----------------------|------------------------------|-----------------|---------------------------|-----------------------------|--------|
| ERR1948631_1.fastq | PASS | FAIL | WARNING | FAIL | 🔍 |
| clean_ERR1948631_1.fq | WARNING | PASS | WARNING | FAIL | 🔍 |

The FASTQ quality check task is performed by nine analysis modules. The table above provides a quick evaluation of whether the results of each module seem entirely normal (pass), slightly abnormal (warning) or very unusual (fail). Note that these evaluations must be taken in the context of what is expected from the library. For example, some experiments may be expected to produce libraries which are biased in particular ways. Therefore, the summary evaluations should be treated as pointers that guide the preprocessing of the libraries.



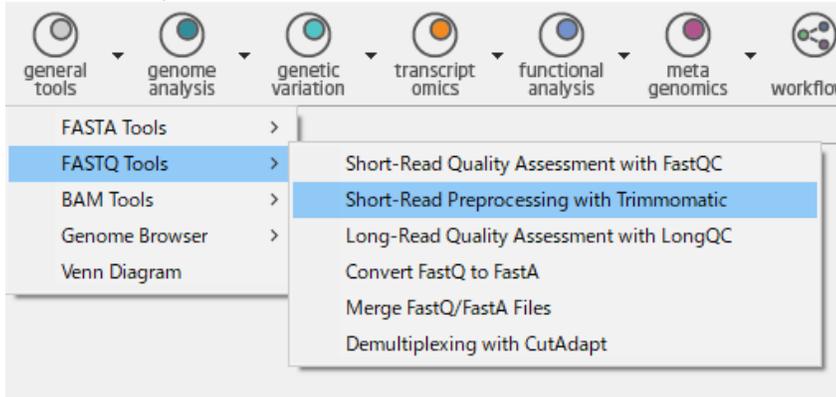
✓ 解析が終了するとレポートが作成

正常 (PASS)
 わずかに異常 (WARNING)
 異常 (FAIL)

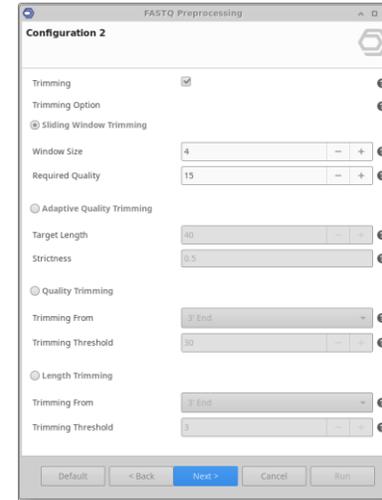
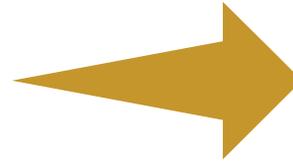
シーケンスデータの品質をすばやく評価



レポートのアイコンをクリック→さらに詳細な結果を見ることが可能



QC後必要に応じてトリミングを行う

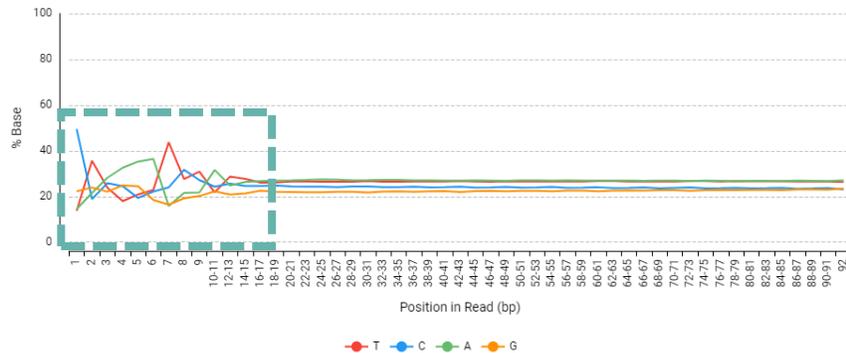


様々なトリミング方法
トリミング後再度QCを行う

【補足】必ずしもすべてのQC項目を正常（PASS）にする必要はありません。

| Name | Per Base Sequence Quality | Per Sequence Quality Scores | Per Base Sequence Content |
|-----------------------|---------------------------|-----------------------------|---------------------------|
| SRR6312175_1.fastq.gz | PASS | PASS | FAIL |
| SRR6312174_1.fastq.gz | PASS | PASS | FAIL |
| SRR6312174_2.fastq.gz | PASS | PASS | FAIL |

Per Base Sequence Content [SRR6312174_1.fastq.gz]

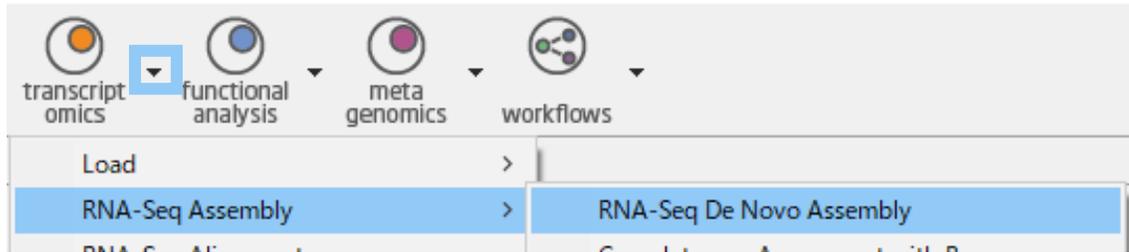


FalseやWarningと判定されていても、それを説明できる理由があれば、適切なリードであると判断することが多いです。

例)

左の結果ではPer base Sequence ContentがFalseですが、こちらはRNA-Seqデータであれば想定される形状に見られます。

RNA-seqライブラリー調製におけるプライミングはそれほどランダムではないので、最初の塩基の塩基含量のシフトが起こることが予想されるためです。



① de novo アセンブリツールをクリックします。

Transcript Moduleを使用します。

Input

The RNA-seq de novo Assembly task consists of reconstructing the transcriptome from RNA sequencing data, assembling short nucleotide sequences into longer ones without the use of a reference genome. This functionality is based on Trinity, a De Bruijn assembler software.

Note: This tool makes use of free cloud computation resources. This is an introductory offer and may change in a future release depending on the overall resource consumption of this feature.

Sequencing Data: Paired-End Reads

Sequencing Format: FASTQ

Input Reads: 4 Files

- /home/marta/Downloads/de_novo_assembly/SRR937568_2.fastq.gz
- /home/marta/Downloads/de_novo_assembly/SRR937568_1.fastq.gz
- /home/marta/Downloads/de_novo_assembly/SRR937564_2.fastq.gz
- /home/marta/Downloads/de_novo_assembly/SRR937564_1.fastq.gz

Paired-End Configuration

Define the pattern to distinguish upstream files from downstream files. The pattern is searched right before the file extension, and the rest of the name should be the same for both files of each sample.

Upstream Files Pattern: _1

Downstream Files Pattern: _2

Default < Back Next > Cancel Run

ペアエンドリードか
シングルエンドリードか
を指定

② 使用するデータをすべて入力します。

RNA-seqリードの鎖長を定義します。
ライブラリー調整の条件や受託元にご確認ください。

コンティグの最小長

Bowtie2を使用して入力RNA-Seqリードを
トランスクリプトームアセンブリにアラインメントします。

アセンブリ結果からSuperTranscriptsを構築します。

正規化をオフ
通常チェックを入れない

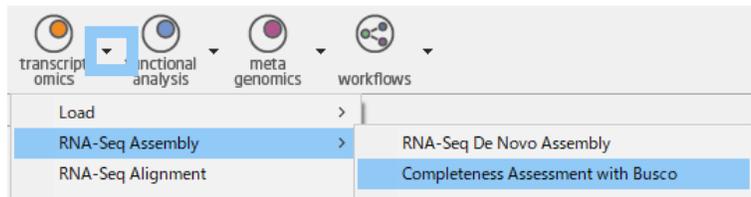
融合転写産物を最小限に抑えることができます
コンパクトな真菌ゲノムでは、強くお勧めします。
注) 負荷の高い操作なので、必要がない限り
使用しないでください。

④ 次の画面のパラメータは、Trinityツールに詳しい
ユーザー向けの設定を目的としています。
(セミナーではデフォルトの設定で実行)

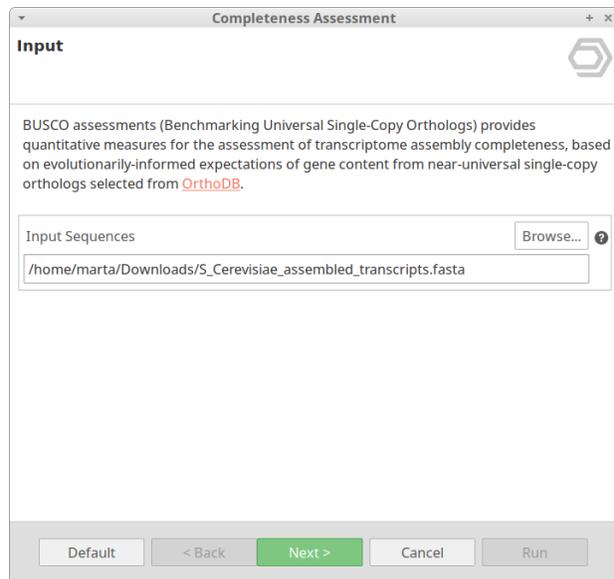
③ 必要に応じてパラメータ設定を変更します。

(セミナーで使用したサンプルは真菌のためPaired-End Configurationのみチェックをいれ後はデフォルトの設定で実行)

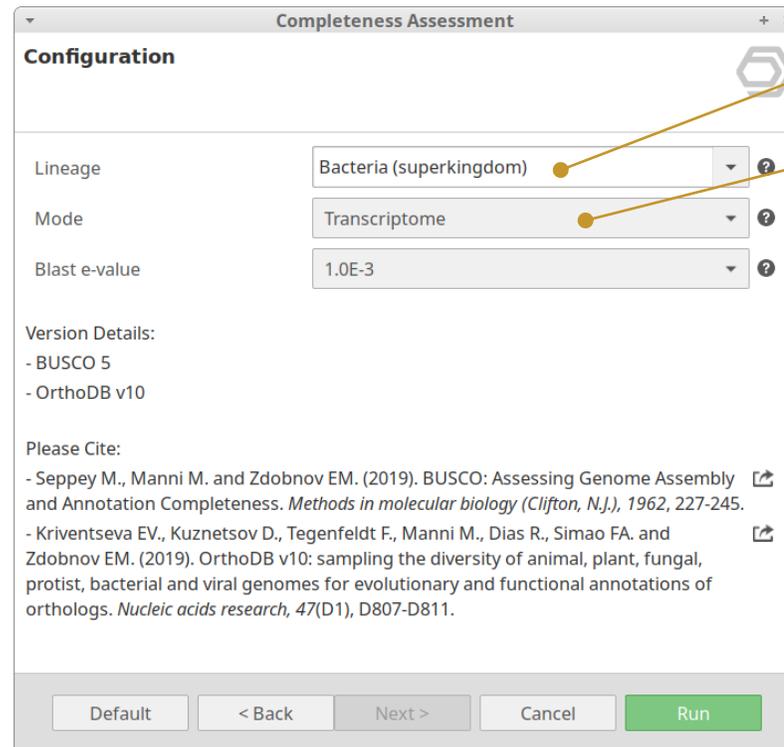
結果のFASTAファイルに加えて、レポートが生成されます。
このFASTAファイルをしようしてコンティグの完全性を評価します。



①完全性評価ツールをクリックします。



②de novo アセンブリ結果のFASTAファイルを入力します。



● 評価する種に応じて、適切な系統を選択します。

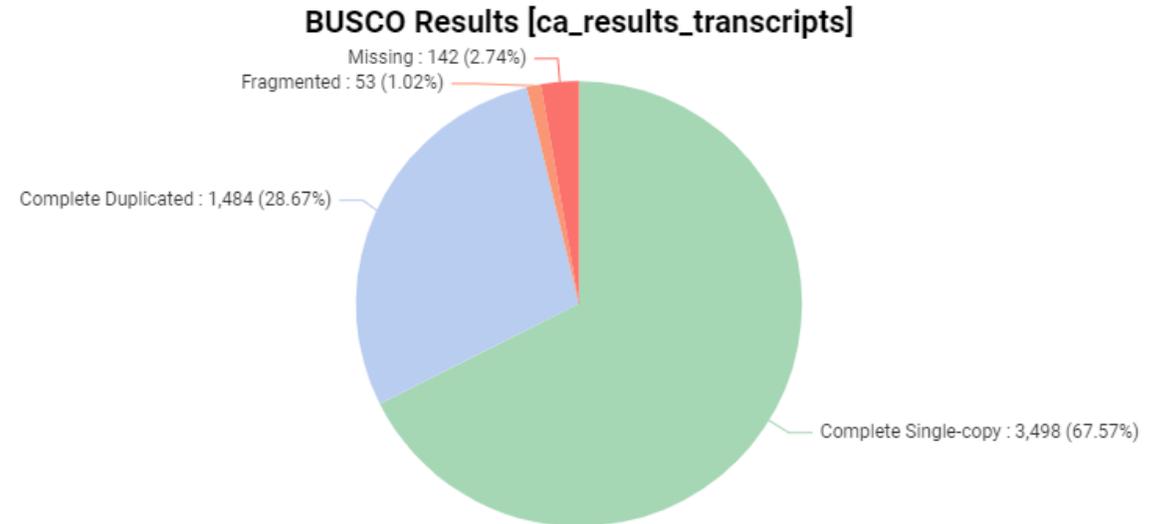
● Transcriptomeを選択します。

③データに合った設定を行います。

BUSCO Results

96.23% of the BUSCO groups have complete gene representation (single-copy or duplicated), whi

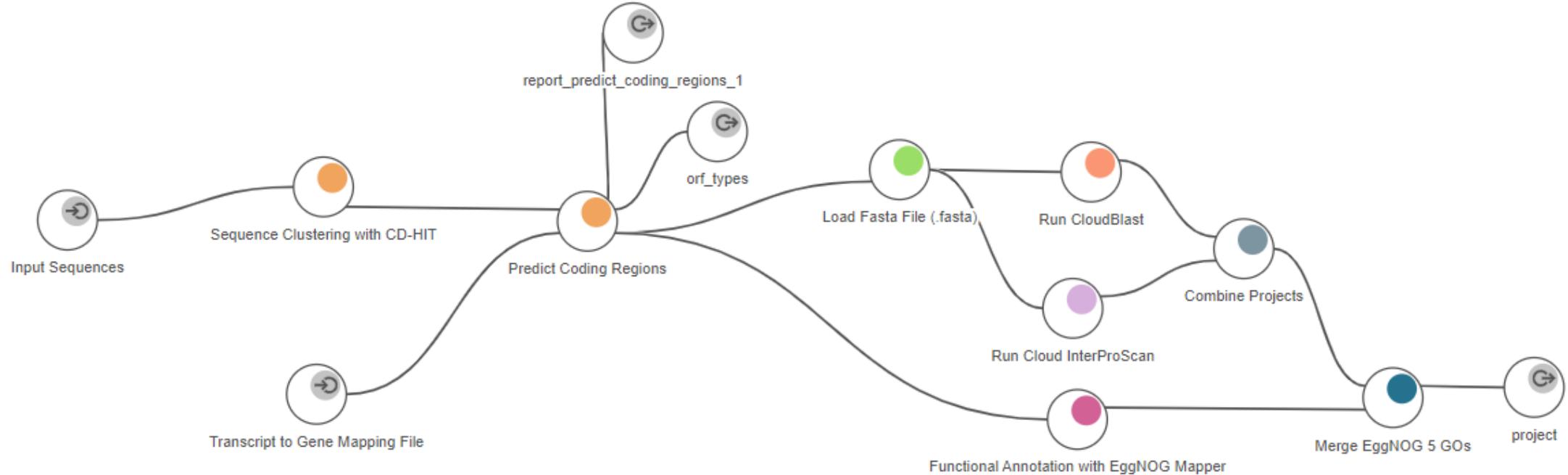
| BUSCO Notation | BUSCOs | BUSCO IDs | Sequences | Sequence IDs |
|----------------------|----------------|-----------|-----------------|--------------|
| Complete Single-copy | 3,498 / 67.57% | Id | 3,359 / 10.4% | Id |
| Complete Duplicated | 1,484 / 28.67% | Id | 3,405 / 10.54% | Id |
| Fragmented | 53 / 1.02% | Id | 53 / 0.16% | Id |
| Missing | 142 / 2.74% | Id | 25,630 / 79.34% | Id |



レポートやチャートからコンティグの完全性を簡単に評価できます。

上記の例ではアセンブルされた配列中のおおよそ96%がCore gene setに含まれていたため、下流分析に使用できると判断することができます。

残りの解析についてはワークフローを使用します。



ワークフローでは各ツールがアイコンで表現されています。
赤枠のアイコンは設定が不十分でワークフローを開始できません。
またアイコンの設定はすべてデフォルトの設定になっているので、各ツールの設定を確認・修正が必要です。
設定はアイコン上で右クリック、もしくはダブルクリックで行います。



の設定画面を開くと次のような画面が現れます。

Sequence Clustering with CD-HIT

Input

CD-HIT is a widely used program for clustering biological sequences to reduce sequence redundancy and improve the performance of other sequence analyses. CD-HIT-EST clusters a nucleotide dataset into clusters that meet a user-defined similarity threshold, usually a sequence identity.

Input Sequences

Default < Back **Next >** Cancel Run

Sequence Clustering with CD-HIT

Configuration 1. Algorithm Options

Sequence Identity Type: Global

Sequence Identity Threshold: 0.95

Band Width: 20

Word Length: 10

Length Cutoff: 10

Length Difference Cutoff: 0

Accurate Mode:

Comparing Both Strands:

Default < Back **Next >** Cancel Run

Sequence Clustering with CD-HIT

Configuration 2. Alignment Coverage Options

Adjust Longer Sequence Coverage:

Longer Sequence Coverage: 0.9

Adjust Shorter Sequence Coverage:

Shorter Sequence Coverage: 0.9

Longer Sequence Unmatched %: 1

Shorter Sequence Unmatched %: 1

Alignment Position Constraints:

Default < Back **Next >** Cancel Run

① de novo アセンブリ結果のFASTAファイルを入力します。

② 任意で設定を行います。

(セミナーではデフォルトの設定で実行)



の設定画面を開くと次のような画面が現れます。

de novo アセンブリの工程で作成された Transcript to Gene Mapping File テキストデータを入力します。

Pfam 検索により、既知のタンパク質との相同性を持つ ORF を特定します。有効にすることをお勧めいたしますが、実行時間が大幅に増加することに注意してください。

①必要に応じてパラメータ設定を変更します。

(セミナーでは Transcript to Gene Mapping File の入力と Pfam 検索のみデフォルト値から変更)

②コーディング領域の配置の可能性を予測するための設定を任意で変更します。

(セミナーでは single Best Only にチェックを入れ
それ以外はデフォルトの設定で実行)



の設定画面を開くと次のような画面が現れます。

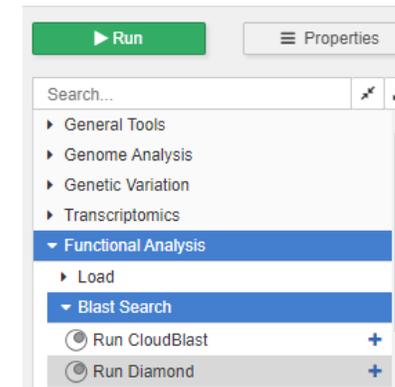
タンパク質配列のためblastp(あるいはblastp-fast)を指定します。

任意のデータベースを指定します。

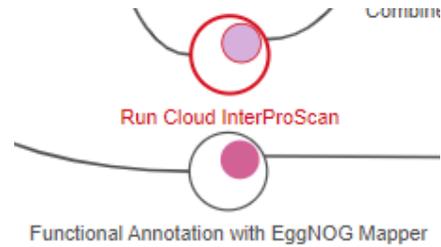
分類を選択することによってより高速のBlast結果が得られ、Cloud Computation Unit*の消費量が少なくなり、機能アノテーションの具体性が高まります。

*Computation Unitsは、CloudBlast解析を行うごとに消費されます。(InterProScanのCloud解析でも消費)すべてのUnitsを使いってしまった場合は、追加で 6 million Computation Units を購入することができます。

【補足】
Cloud BlastをDiamondに変更できます。

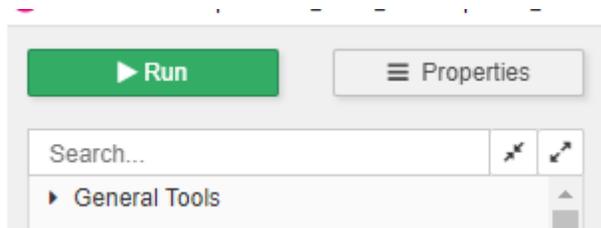


ワークフロー左側のメニューからRun Diamondの「+」をクリックすると解析アイコンが出現します。後はCloud Blastを削除しこのアイコンに置き換えます。



の設定も任意に変更可能 (セミナーではデフォルトの設定で実行)

*InterProScanは保存先を指定するだけで赤枠のアイコンから通常のアイコンに変わります。



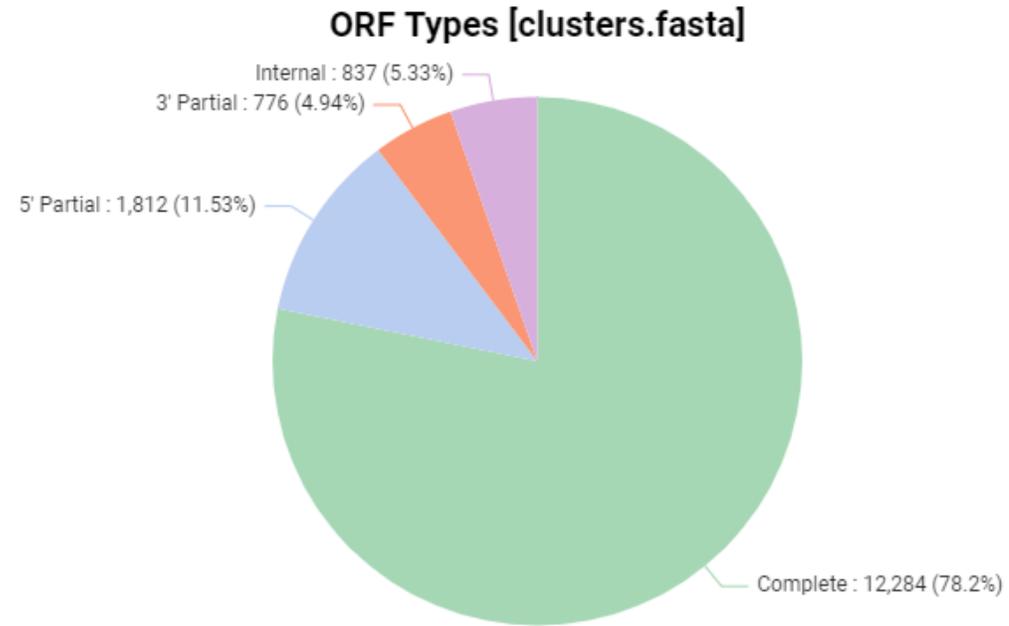
全ての設定が完了したら、Runボタンをクリックすることで一連の解析がスタートします。

コーディング領域予測の結果

TransDecoder Results

| ORF Type | Sequences | |
|------------|-----------|----|
| Complete | 12,284 | ld |
| 5' Partial | 1,812 | ld |
| 3' Partial | 776 | ld |
| Internal | 837 | ld |

Longest ORF per gene: lv



レポートやチャートから予測されたコーディング領域属性を簡単に評価できます。

Complete: 開始コドンと停止コドンが含まれています。

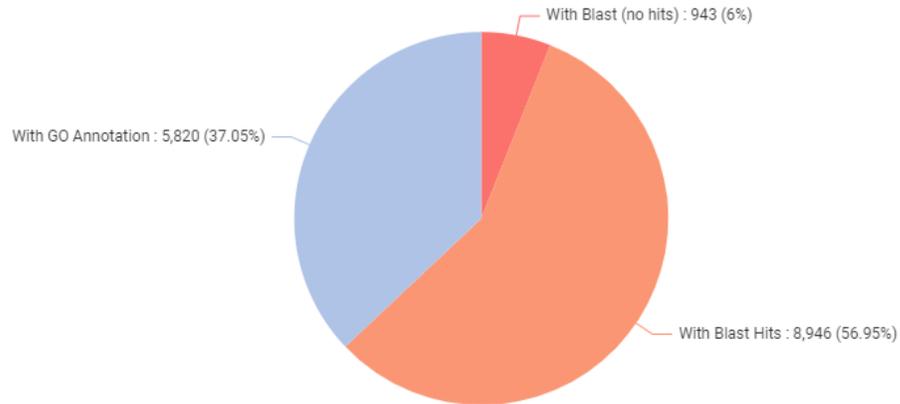
5' partial: 開始コドンが欠落しており、おそらく N 末端の一部が欠落しています。

3' partial: 終止コドンが欠落しており、おそらく C 末端の一部が欠落しています。

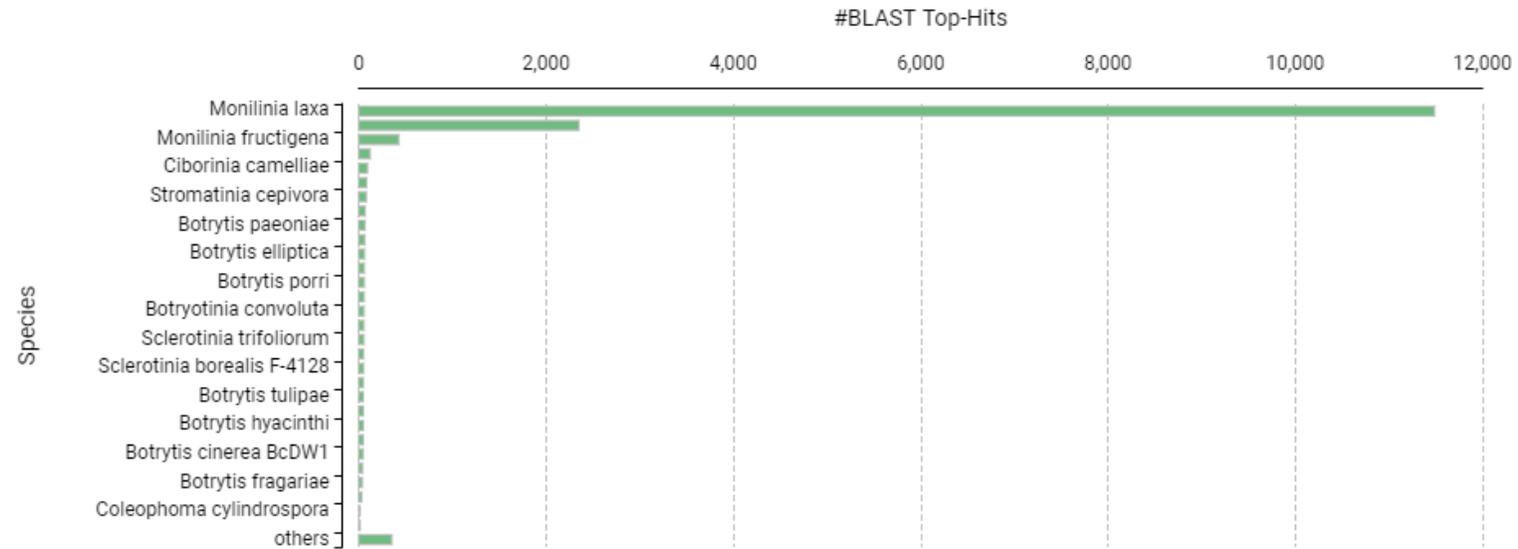
Internal: 5' と 3' の両方とも部分的です。

Blat、InterProScan、EggNog Mapperの結果

Data Distribution Pie Chart [project]

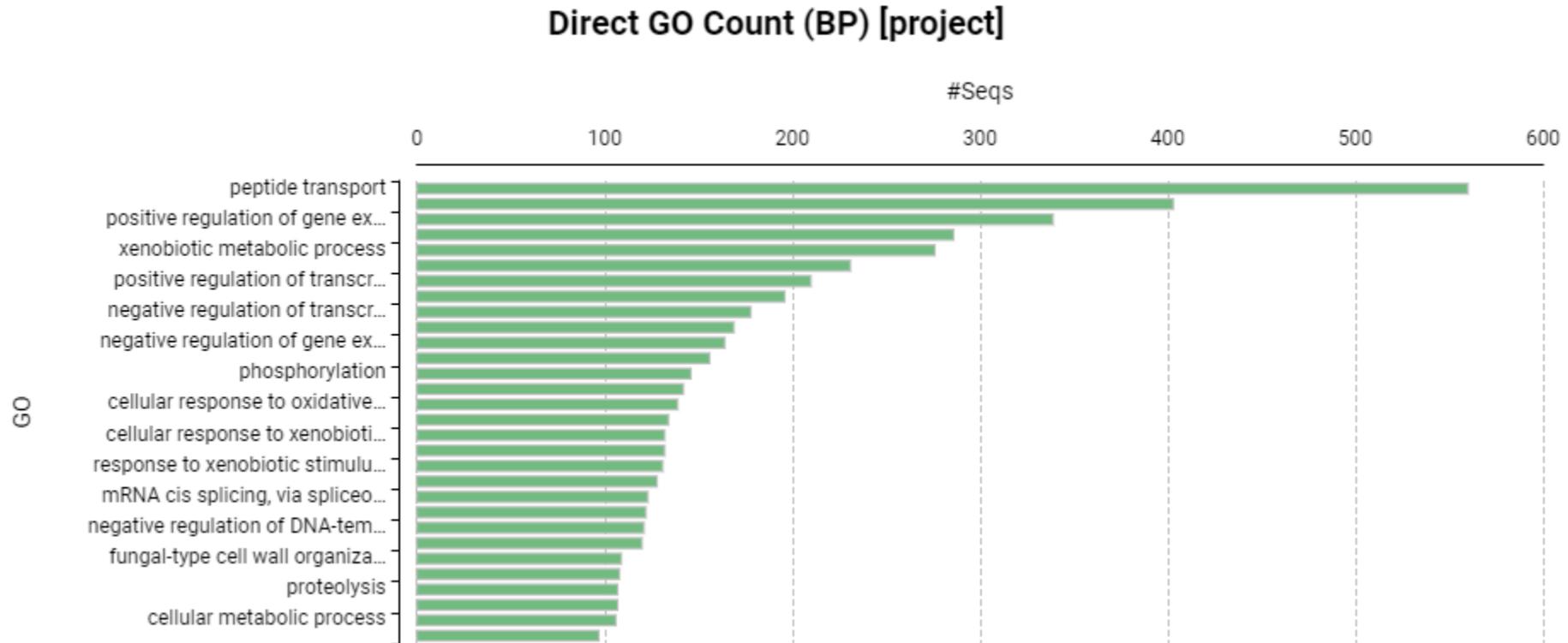


Top-Hit Species Distribution [project]



どのぐらいのトランスクリプトーム配列に特徴がつけられたか
どの生物種のアノテーションが付与されたかがグラフで表現できます。

Blat、InterProScan、EggNog Mapperの結果



データセット内で最も頻繁に使用される GO Termを表示します(上図はbiological processカテゴリ)。

お問い合わせ先：フィルジェン株式会社

TEL 052-624-4388 (9:00～17 : 00)

FAX 052-624-4389

E-mail: biosupport@filgen.jp