

## 次世代シーケンサーデータの解析手法 第 12 回 Galaxy : ヒストリーとワークフロー

寺田 朋子<sup>1</sup>、大田 達郎<sup>2\*</sup>、清水 謙多郎<sup>1,3</sup>、門田 幸二<sup>1,3\*</sup>

<sup>1</sup> 東京大学 大学院農学生命科学研究科

<sup>2</sup> 情報・システム研究機構 データサイエンス共同利用基盤施設  
ライフサイエンス統合データベースセンター

<sup>3</sup> 東京大学 微生物科学イノベーション連携研究機構

Galaxy は、ウェブブラウザ上でマウスを操作して行う GUI ベースの (Linux コマンドを覚える必要のない) データ解析環境である。今回は、Galaxy の概要と公共サーバ Galaxy Main の基本的な利用法について述べた。第 12 回は、前回 Galaxy Main 上で行った解析結果 (ヒストリー) の他の研究者との共有や、以前行った解析手順 (ワークフロー) を他のデータに適用する手段を中心に述べる。また、ヒストリー間のデータのコピーや、自身の PC を経由しない Galaxy Main へのデータの取り込みなど、関連した便利な機能についても紹介する。今回の内容は、ウェブブラウザの違いに起因する不具合を避けるため、Google Chrome または Firefox (Internet Explorer は非推奨) を用いてほしい。ウェブサイト (R で) 塩基配列解析 (URL: [http://www.iu.a.u-tokyo.ac.jp/~kadota/r\\_seq.html](http://www.iu.a.u-tokyo.ac.jp/~kadota/r_seq.html)) 中に本連載をまとめた項目 (URL: [http://www.iu.a.u-tokyo.ac.jp/~kadota/r\\_seq.html#about\\_book\\_JSLAB](http://www.iu.a.u-tokyo.ac.jp/~kadota/r_seq.html#about_book_JSLAB)) が存在する。ウェブ資料 (以下、W) や関連ウェブサイトなどを効率的に活用してほしい。

Key words : NGS, Galaxy, workflow, history

### はじめに

本連載では、これまで主にウェブブラウザ Internet Explorer (IE) を用いて動作確認を行ってきた。しかしながら、今回初めて Galaxy<sup>1)</sup> の一部機能 (ヒストリーの一覧表示やワークフロー) の実行時に、IE の使用に起因する不具合に遭遇した。2018 年 6 月現在の最新 OS である Windows 10 は、推奨ウェブブラウザが Microsoft Edge である。しかしながら、使い慣れた IE を継続的に利用している Windows ユーザはまだ相当いると思われる。この問題は、おそらく IE ユーザに限定されるものの、原因特

定に要する労力および判明時の落胆が非常に大きいため注意してほしい。我々は今回の動作確認を Chrome で行ったが、Firefox でも問題ないと思われる [W1]。

今回も、前回 (第 11 回) に引き続いて Galaxy Main 上で作業を行う。必ずしも予め前回の内容を一通り行っておく必要はない。次項で述べるヒストリーの共有によって、前回までの解析結果 (ヒストリー) を共有することができるからである。但し、最低限第 11 回の W4 を参考にして Galaxy Main のユーザ登録を済ませておくことを推奨する。Galaxy Main にログインした状態で作業を行うことで、ヒストリーやワークフローに関する今回の内容を十分に活かすことができるからである。尚、ここでは Galaxy Main にユーザ登録済のメインユーザを kadota@iu.a.u-tokyo.ac.jp (以下、kadota\_registered) とする [W2]。また、ヒストリーを共有する登録済のユーザとして tomoko.terada@iu.a.u-tokyo.ac.jp (以下、terada\_registered) を、そして未登録ユーザを koji.kadota@gmail.com (以下、

\*To whom correspondence should be addressed.

Phone : +81-3-5841-2395

Fax : +81-3-5841-1136

E-mail : tohta@dbcls.rois.ac.jp

kadota@bi.a.u-tokyo.ac.jp

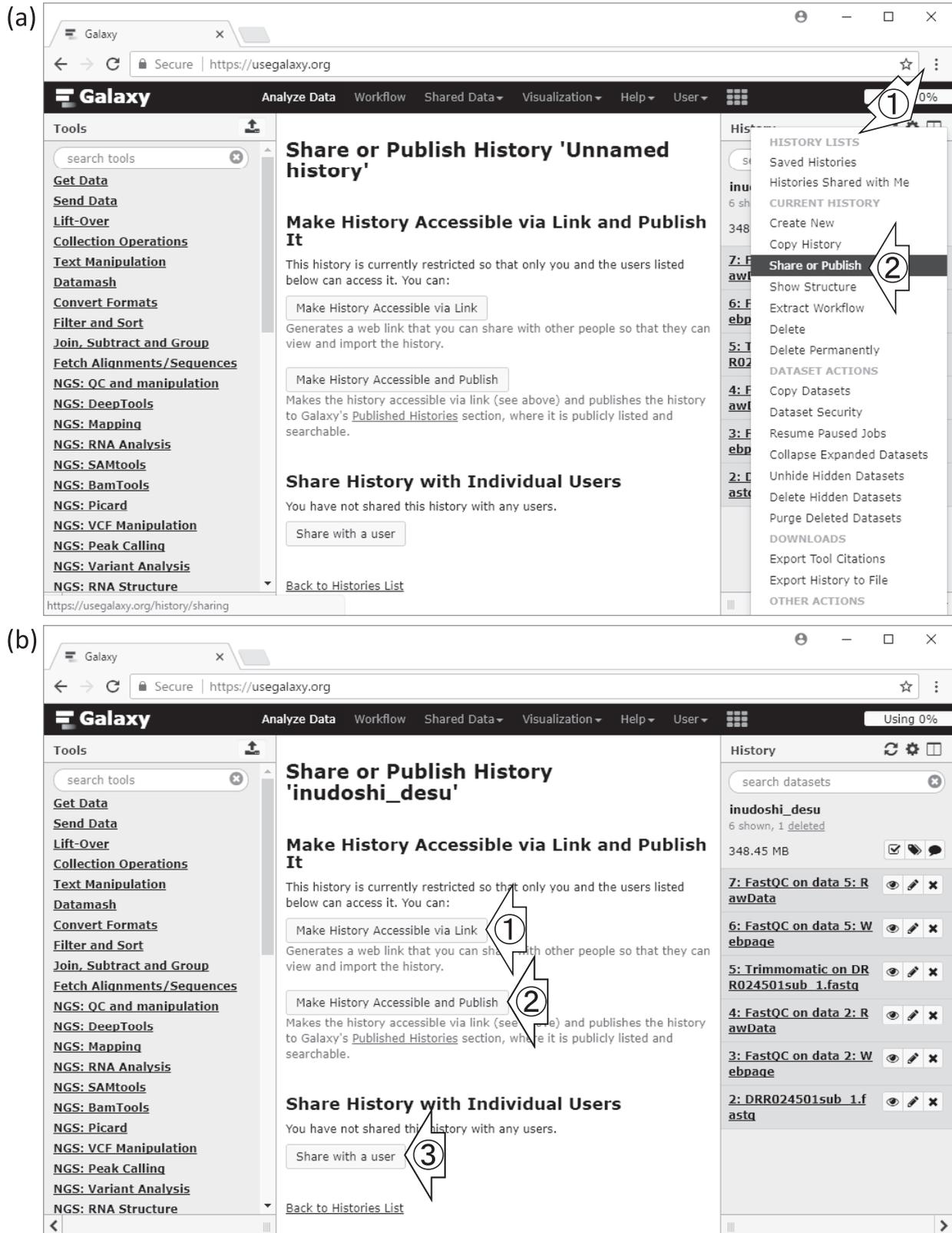


図 1. ヒストリーの共有 (その 1)

(a) は W5-3 と、そして (b) は W5-4 と基本的に同じである。(a) でヒストリー名を “inudoshi\_desu” に変更後、① History options - ② Share or Publish を実行すると、(b) のように変更が反映される。(b) で① Make History Accessible via Link をクリックした結果が図 2a である。

kadota\_unregistered)として話を進める。

## 履歴の共有

前回までに kadota\_registered が一通り行った作業は、Illumina MiSeq を用いて得られた 30 万リードからなる gzip 圧縮 FASTQ ファイル (DRR024501sub\_1.fastq.gz) のアップロード、FastQC の実行、Trimmomatic の実行、そして Trimmomatic 実行後のファイルに対する FastQC の実行であった [W3]。Galaxy GUI 画面の右側にある履歴パネル上で「History options」-「Share or Publish」を選ぶと、中央のパネル上に 3 つの履歴共有手段が提示される [W4-3] (図 1)。

- ① Make History Accessible via Link  
リンク先 URL を知っているヒトのみで履歴を共有
- ② Make History Accessible and Publish  
Galaxy Main の Published Histories というサイト上で公開
- ③ Share with a user  
特定のユーザのみと共有

①と③の違いが分かりづらいかもしれないが、①のほうはリンク先の URL を知っていれば誰でも履歴情報にアクセスできるため、③よりも制限が緩いということになる。しかしながら、その URL を知っているのは教えられたヒトのみであるため、通常利用の局面においては実質的に同じであろう。②を選択すると、Galaxy Main の Published Histories というサイトで公開される。このサイトでは様々な履歴名のもが見られるが、たとえ同じ研究グループ内のヒトであったとしても、履歴内部の理解は簡単ではないだろう [W4-5]。尚、この履歴名 (デフォルトは“Unnamed history”) は、右側の履歴パネル上で任意の名前に変更することができる [W5-2]。変更後の履歴名“inudoshi\_desu”で、改めて「History options」-「Share or Publish」を選ぶと、中央のパネルのほうも反映される [W5-3]。

ここでは、連載第 11 回で行った kadota\_registered の解析結果 (履歴) を、未登録ユーザである kadota\_unregistered と共有する①のやり方を示す。図 1b において①をクリックすると、履歴 inudoshi\_desu のリンク先 URL が現れる [W6-1] (図 2a)。履歴情報を提供する側である kadota\_registered の仕事は、基本的にこの URL 情報を共有したいヒトにメール送信するだけである [W6-2]。図 2b は、情報を提供される側の kadota\_unregistered がこの URL にアクセスした結果である [W6-4]。提供する側 (kadota\_registered) の履歴パネル情報に相当するものが、左のほうに見えていること

がわかる。提供する側と見栄えが異なるのは、編集権限を持たない閲覧専用のような状態でこの履歴情報を眺めているからである。この状態を解除するには、Import history ボタンを押して履歴 inudoshi\_desu を取り込めばよい (インポートすればよい) [W7-1]。これは、Galaxy Main にログインしていない kadota\_unregistered の Galaxy 環境に、履歴情報 inudoshi\_desu をコピーしていることに相当する。

提供された側 (kadota\_unregistered) のインポート後の画面は、提供した側 (kadota\_registered) の Galaxy Main の通常画面と同じである [W7-2]。共有された履歴情報をインポートすることで、基本的に左側のツール選択パネルから自由にプログラムを選択して、独立に解析を行うことができるようになる。もちろん Galaxy Main にログインしていない状態のため、インポートした履歴やその後独立して行った解析結果は、ウェブブラウザ上の開いているタブまたはブラウザそのものを閉じれば消える [W7-3]。これをどう捉えるかはヒトそれぞれであるが、例えば次のような利用法が考えられる。研究グループ内に Galaxy を使えるヒト (例: kadota\_registered) がおり、一通りのデータ解析結果が Galaxy Main 上に存在するとしよう。データ解析担当者が、Galaxy Main 上でのデータ解析結果 (の履歴の URL 情報) をグループ内全員に知らせることで、グループ内の他の研究者 (例: kadota\_unregistered) が解析手順にミスがないかなどをチェックすることができる。

必要に応じてインポートを行うことで、例えば Trimmomatic<sup>2)</sup> 実行時にデフォルトの平均クオリティスコアを 20 ではなく 30 にすればどうなっていたかなどを独自に調べることもできる [W9-6]。但し、Galaxy Main にログインしていない状態では機能が制限されるため、kadota\_registered の履歴 inudoshi\_desu に対して kadota\_unregistered が追加作業を行った履歴を、さらに他者と共有することはできない [W10-1]。もちろん、この追加作業情報を保存することはできる。追加作業を行ったブラウザ上で、アカウントがなければ新規作成するなどして Galaxy Main にログインすればよいのである [W11-1]。

## 履歴の操作

図 3a は、kadota\_unregistered が履歴 inudoshi\_desu をインポートし、追加作業を行った (オプションを変更して Trimmomatic および FastQC を実行した) 結果のブラウザ上で、kadota\_registered がログインした直後の画面である [W11-2]。kadota\_registered が前回のログイン時に行った作業は、履歴 inudoshi\_desu の作成までであった [W6-1]。kadota\_registered でログイン後の履歴パネルは、ログイン前に kadota\_unregistered

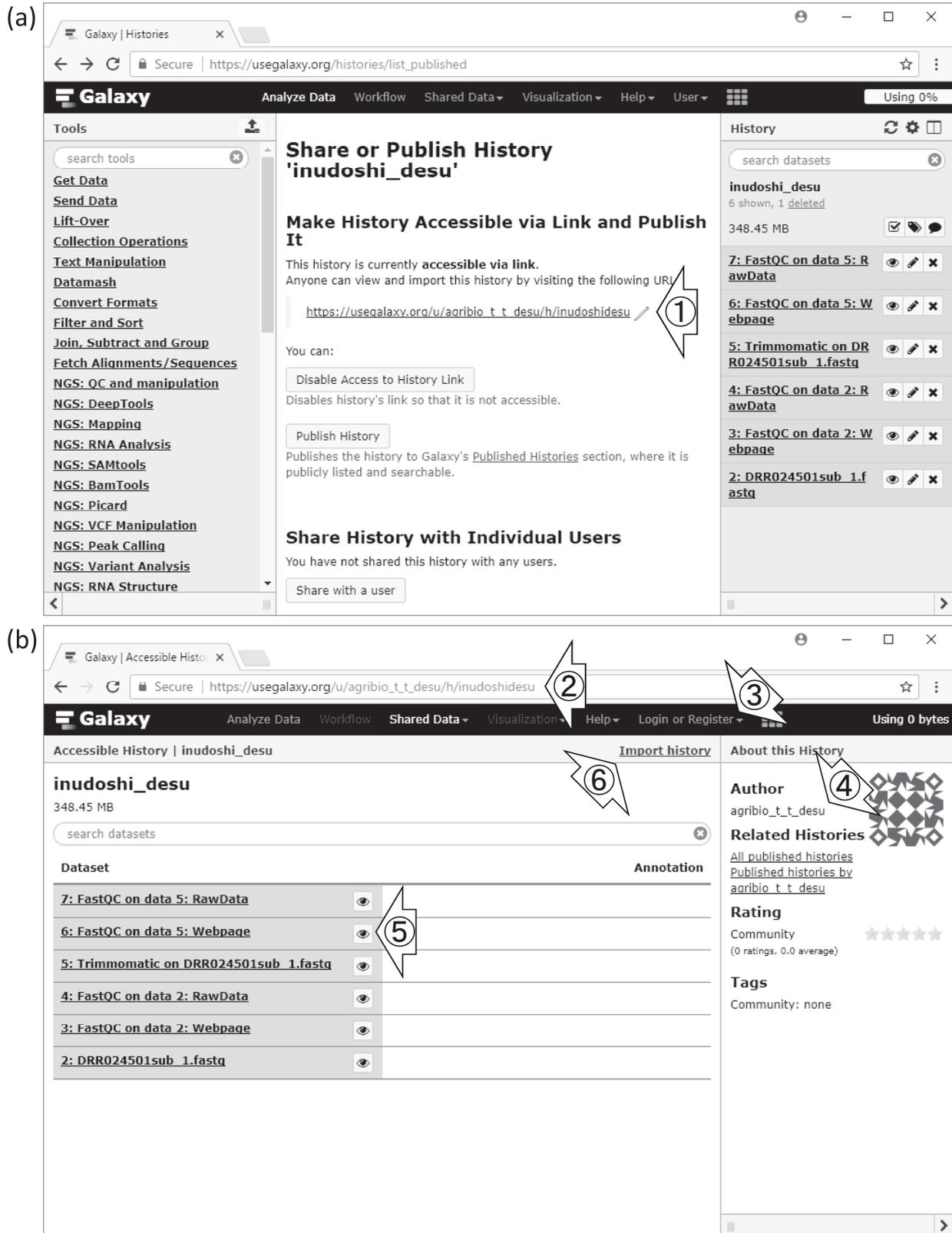


図 2. ヒストリーの共有 (その 2)

(a) は W6-1 と、そして (b) は W6-4 から 6-6 と同じ。(a) はヒストリー情報を提供する側の画面、(b) は提供された側がリンク先にアクセスした結果の画面。①はこのヒストリーの URL 情報で、②と同じ。③確かにログインしていないことがわかる。④はヒストリー提供者 kadota\_registered の Public name。⑤をクリックすると FastQC 実行結果が見られる。⑥ Import history でヒストリー情報を取り込むこともできる。

が追加作業を行った W9-4 までと同じである。kadota\_unregistered による履歴 inudoshi\_desu のインポート前後の違いからもわかるように [W7]、kadota\_registered のログインによって追加作業情報を取り込んでいることを意味する。

もちろん、図 3a の①履歴情報は、W6-1 までで作成した履歴 inudoshi\_desu とは別物である。履歴パネルの② View all histories を実行すると、①現在の履歴情報 (Current History) が左側に、そして③ W6-1 までで作成した履歴 inudoshi\_desu が右側に表示されていることがわかる (図 3b)。履歴の削除・切替・新規作成・反映などの基本操作 [W12] 以外の有効なテクニックとしては、履歴間でのデータコピーが挙げられる [W13-3]。他の履歴から現在の履歴への一方向のコピーしかできないが、上記の基本操作 (履歴の切替と反映) を組み合わせれば問題ない。具体的な利用法としては、例えば新規履歴作成後に手元のデータを改めてアップロードすることなく、既存の履歴内にあるデータをコピーして任意の解析をスタートさせることなどが挙げられる [W13-7]。

## ワークフローの作成

ワークフローは、以前行った解析手順を別のデータに対して実行したい場合に有効な手段である。履歴 inudoshi\_desu は、Illumina MiSeq を用いて得られた 30 万リードからなる paired-end の forward 側データ (DRR024501sub\_1.fastq.gz) に対して解析を行ったものである [W14-1]。ここでは、inudoshi\_desu をベースにワークフローを作成し、第 6 回<sup>3)</sup>の W3-2 で作成した reverse 側のファイル (DRR024501sub\_2.fastq.gz) に対して独立に行うやり方を説明する。

ワークフローは、作業の流れを示したものであるため<sup>4)</sup>、白紙のノートに利用するプログラムを 1 つずつ選択して手順書を新規作成するように、ワークフローを 1 から作成することはもちろんできる。しかしながら、作成したいワークフローと似た既存の履歴を利用するのがおそらく一般的である。ここでは、「FastQC → Trimmomatic → FastQC」の 3 ステップからなるオリジナルの inudoshi\_desu に対して、1 ステップ (平均クオリティスコアの閾値を 30 に変更した Trimmomatic) 追加して実行したもの [W13-3] をベースとして用いるべく Current History にした [W14-2]。

現在履歴パネル上に表示されているものをワークフローのベースとして利用するには、「History options」-「Extract Workflow」を選択すればよい [W14-3] (図 4a)。中央パネルに右側の履歴パネルに対応したワークフローが表示されるので、ワークフロー名などを任意に変更して編集する [W14-4]。この履歴に対応したワー

クフロー (作業の流れ) は、「FastQC → Trimmomatic (default で実行) → FastQC → Trimmomatic (オプション変更して実行)」の 4 ステップからなる。ステップごとに「Include “プログラム名” in workflow」のチェックのオンオフを行えるので [W14-8]、ここでは 4 ステップ目の Trimmomatic のチェックをオフにした、3 ステップからなるワークフローを作成した [W14-9] (図 4b)。左側のツール選択パネルから利用できるようにしておくと、使い勝手がよいだろう [W15-4]。

## ワークフローの実行

reverse 側のファイル (DRR024501sub\_2.fastq.gz) を入力として、作成したワークフローを実行する手順を説明する。入力ファイルとして用いることができるのは履歴パネル上で見られるものに限定されるため、まずはその作業を行う [W16-2]。ローカル環境 (自分の PC 内) にあるファイルをアップロードするやり方については、第 11 回の W5-3 で示した。ここでは、指定した URL 上にあるファイルを Galaxy 上にダウンロードするやり方を示す。アップロードとダウンロードの言葉の違いは本質的ではなく、Galaxy への他の便利なデータ取り込み手段を知るのが目的である。URL を指定するやり方の一番のメリットは、自分のローカル環境にデータを一旦保存する必要がない点である [W17-8]。W17-5 で示した URL は所属機関のものであるが、第 6 回 W3-1 で用いた DDBJ SRA<sup>5)</sup> の ftp サイトの URL など置き換えて考えれば、その便利さが想像できるかもしれない。

図 5 は、ワークフローの (a) 実行直前と (b) 実行後の画面である。新規履歴を作成し [W17-1]、reverse 側のファイル (DRR024501sub\_2.fastq.gz) のみを Galaxy 上にダウンロードした状態にしておくことで [W17-10]、ワークフローの実行をスムーズに行うことができる。計 3 ステップからなるワークフローのうち、Step1 と Step2 は独立に実行可能である一方、Step3 (2 回目の FastQC) は Step2 (Trimmomatic) の実行結果を入力データとして利用する関係となっている。このような Step2 と Step3 の依存関係についても、ワークフロー中に正しく記載されているので心配ない [W18-6]。ワークフロー実行ボタンを押すと、3 ステップ分のプログラムが文字通り流れ作業的に実行される。このときは Step2 が Step1 よりも先に終了したが、上記の理由 (独立に実行可能) により全く問題ない [W18-10]。このワークフローの場合は、Step3 が Step2 よりも後に実行されるはずであり、実際そうになっていることがわかる [W18-12]。

## ワークフローの編集と共有

ワークフロー内のステップ間の依存関係については、

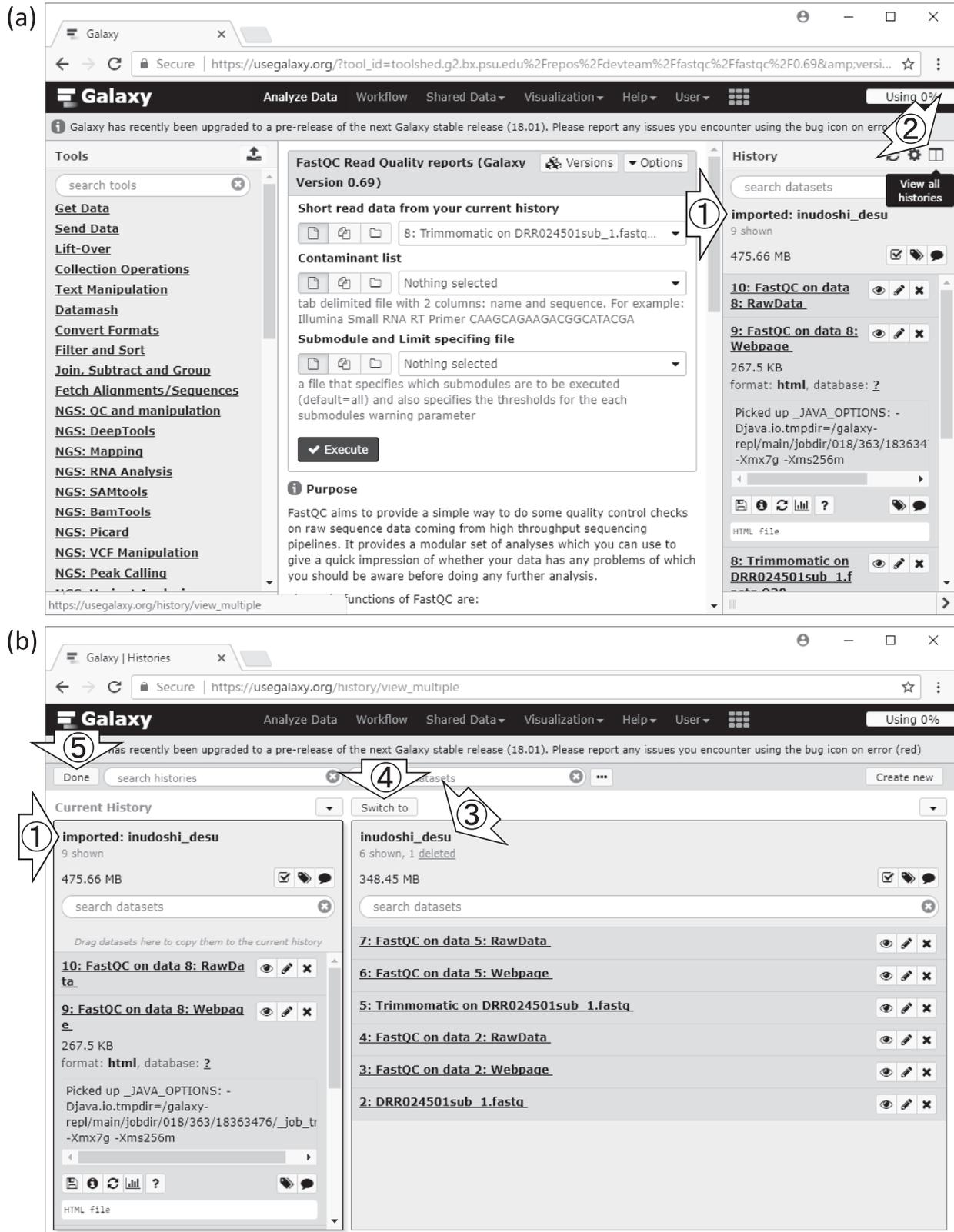


図 3. ヒストリーの操作

(a) は W11-2 と、そして (b) は W11-4 と同じ。①が現在のヒストリー情報。② View all histories をクリックすると、(b) の画面に切り替わり、kadota\_registered の全ヒストリー情報が見られる。③ W6-1 までで作成したヒストリー inudoshi\_desu を現在のヒストリーに切り替えたい場合は、④ Switch to、⑤ Done をクリックすればよい。



図 4. ワークフローの作成

(a) は W14-3 と、そして (b) は W14-9 と同じ。(a) 作成したいワークフローに近いヒストリーをひな形として利用すべく、Current History にした状態で① History options ➡ ② Extract Workflow。(b) 中央パネルの編集画面上で利用するプログラムを選択したのち、③ Create Workflow すれば④の名前のワークフローが作成される。



図 5. ワークフローの実行

(a) は W18-7 と、そして (b) は W18-12 と同じ。(a) 新規ヒストリーを作成し (W17-1)、① reverse 側のファイルを Galaxy 上に取り込んで (W17-10)、中央パネル上で② Step3 の入力ファイルが③ Step2 の実行結果となっていることを確認したのち、④ワークフローを実行。(b) 無事終了した状態。

ワークフローの編集画面を眺めると理解しやすいかもしれない [W19-1]。このワークフローの場合、Step2 と Step3 が線で結ばれており、Step2 の出力が Step3 の入力として利用されていることが視覚的にわかる [W19-7]。ワークフロー内で用いられるプログラム内部のオプションについても、ステップごとに変更可能である [W19-8]。また、このワークフローが内部的に FastQC や Trimmomatic を含んでいることがわかるように、任意のキーワードを属性 (Attributes) 情報としてタグ付けすることもできる [W19-9]。

ワークフローの共有手段は、履歴の共有手段 (図 1b) と基本的に同じである。ここでは、特に隠す必要もないので Make Workflow Accessible and Publish を選択して公開した [W20-1]。公開完了後は Galaxy Main の Published Workflows というサイトで公開される。当該ワークフローの URL 情報を他の共同研究者に知らせてもよいし、ワークフロー名の一部の文字でキーワード検索してもらってもよいだろう [W20-3]。

最後に、他の Galaxy Main 登録済ユーザへの公開したワークフロー情報の告知、およびワークフロー利用の実例を示す。ここでは、公開済ワークフローと動作確認用 FASTQ ファイルの URL 情報を `terada_registered` にメールし、`terada_registered` がログインするところから示した [W21-2]。Galaxy Main 登録済ユーザの場合は、一般に履歴パネルに過去の実行結果が残っている状態 (`terada_registered` の場合は履歴 “Trim”) からスタートする [W21-3]。ワークフロー情報の URL にアクセスし、Import workflow ボタンを押すことでワークフロー

が利用可能となる [W21-6]。

動作確認用データのダウンロードが完了すれば、下準備は完了である [W22-1]。ここでは新規履歴 “inudoshi\_desu reverse\_data” を作成して実行結果を保存すべく、履歴オプションを変更して実行した (Send results to a new history を Yes にして Run workflow ボタンを押した)。これは、現在履歴パネル上に表示されている “Trim” には保存されないようにしたことを意味する [W22-8]。正常にワークフローを実行し終えたかどうかを確認するためには、View all histories ボタンを押して、作成した履歴 “inudoshi\_desu reverse\_data” の実行状況を確認する必要がある [W23-2]。W12 でも示したように、現在の履歴 (Current History) に切り替えることで、ワークフロー実行結果を眺めることができる [W23-5]。

## おわりに

今回は、Galaxy の主な特徴である履歴とワークフローについて述べた。例示した内容は非常に簡単なものであったが、実際の解析現場で行われるものも基本的な枠組みは同じである。共有された履歴の内容については、提供側から詳細な説明を受ければよい。自分がワークフローを構築する側になったとしても、公開されているワークフローや自分が以前作成したワークフローの中から目的に近いものを改変して利用すればよい<sup>6)</sup>。また、それが複雑なワークフローの場合でも、step-by-step で組み立てていけばよい。

## 参考文献

- 1) Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, et al. (2016) The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res* 44: W3-W10.
- 2) Bolger AM, Lohse M, Usadel B. (2014) Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30: 2114-2120.
- 3) 谷澤靖洋, 神沼英里, 中村保一, 清水謙多郎, 門田幸二 (2016) 次世代シーケンサーデータの解析手法: 第6回ゲノムアセンブリ. *日本乳酸菌学会誌* 27: 41-52.
- 4) 大田達郎, 寺田朋子, 清水謙多郎, 門田幸二 (2017) 次世代シーケンサーデータの解析手法: 第11回統合データ解析環境 Galaxy. *日本乳酸菌学会誌* 28: 167-175.
- 5) Mashima J, Kodama Y, Fujisawa T, Katayama T, Okuda Y, et al. (2017) DNA Data Bank of Japan. *Nucleic Acids Res* 45: D25-D31.
- 6) 孫建強, 湯敏, 西岡輔, 清水謙多郎, 門田幸二 (2014) 次世代シーケンサーデータの解析手法: 第2回 GUI 環境からコマンドライン環境へ. *日本乳酸菌学会誌* 25: 166-174.

## Methods for analyzing next-generation sequencing data XII. Galaxy – Sharing histories and workflows

Tomoko Terada<sup>1</sup>, Tazro Ohta<sup>2</sup>, Kentaro Shimizu<sup>1, 3</sup>,  
and Koji Kadota<sup>1, 3</sup>

<sup>1</sup> *Graduate School of Agricultural and Life Sciences, The University of Tokyo.*

<sup>2</sup> *Database Center for Life Science, Joint Support-Center for Data Science  
Research, Research Organization of Information and Systems.*

<sup>3</sup> *Collaborative Research Institute for Innovative Microbiology,  
The University of Tokyo.*

### Abstract

Galaxy is an integrative data analysis environment run on the web browser which users can use without using Linux command line. The previous article showed an introduction to the Galaxy system and the basic usage of the public Galaxy server “Galaxy Main.” In this article, using the last article’s results, we present the features to share the analysis results (history) with the other users, or ones to apply the analysis procedures (workflow) to the other data. We also show the useful Galaxy’s features such as copying data across histories, or the direct data import from remote servers. We found a compatibility issue on Internet Explorer with the Galaxy system. Thus we recommend using Google Chrome or Firefox to try the procedures we show in this article. Supplementary materials are available online at: [http://www.iu.a.u-tokyo.ac.jp/~kadota/r\\_seq.html#about\\_book\\_JSLAB](http://www.iu.a.u-tokyo.ac.jp/~kadota/r_seq.html#about_book_JSLAB).