

## 異質倍数体の遺伝子発現量の定量化

山田 恵<sup>†1</sup> 清水(稲継)理恵<sup>†2</sup>  
清水 健太郎<sup>†2</sup> 瀬々 潤<sup>†3</sup>

本研究では、異なる種のゲノムを併せ持つ種である異質倍数体の遺伝子発現量を RNA-seq により定量化した。しかし、異質倍数体の親となる種は近縁の種で構成されることが多いため、通常の RNA-seq では親種同士のホモログ遺伝子の配列が非常に近く、いずれの遺伝子が発現したかを見分けることが困難である。本研究では、近縁種からの変異のパターンに着目し、最尤推定法を用いることで、親種を分けた発現量の定量化を可能とした。

## Quantified expression levels of genes in allopolyploid species

MEGUMI YAMADA,<sup>†1</sup> RIE SHIMIZU(INATSUGU),<sup>†2</sup>  
KENTARO SHIMIZU<sup>†2</sup> and JUN SESE<sup>†3</sup>

In this research, We tackled a measurement of the expression levels of genes in allopolyploids. For the measurement, we used an RNA-seq. However, high sequence similarity between genes originating from different parental species prevented us to apply the method to the problem directly. To overcome the difficulty, our method used genome of model species, which is closely related to the target allopolyploid, and aligned RNA-seq sequences to the genome. Based on the mutation distributions of the alignment result, our method determined the originated species of each read using maximum likelihood estimation method. By counting the number of the sequences of each parental species on each gene, our method successfully quantified expression levels of genes in allopolyploid species.

<sup>†1</sup> お茶の水女子大学 (Ochanomizu University)  
<sup>†2</sup> チューリッヒ大学 (University of Zurich)  
<sup>†3</sup> 東京工業大学 (Tokyo Institute of Technology)

### 1. はじめに

異質倍数体は異なる種のゲノムを併せ持つ種で、ゲノム重複に似た状態を有するため、種分化や環境適応のメカニズム発見に重要と考えられる。これらのゲノムが環境に応じて同調して働くか、個々に働くかの調査に向けて、本研究では、由来するゲノムを区別して遺伝子の発現を同定する手法を開発した。異質倍数体は、植物で多く観測されるため、本論では、*A. thaliana* (シロイヌナズナ) の近縁で起こった異質倍数体の解析を考える。

#### 1.1 ホメオログ

種 A, B を親に持つ異質倍数体 C が存在するとき、A, B がそれぞれ遺伝子 a, b を持つとすると、C は、a, b に由来する  $c_a, c_b$  を持つ (図 1)。a, b をホメオログとすると、 $c_a, c_b$  はホメオログという。A, B は近縁の種である場合が多く、ホメオログは非常に配列相同性が高い。

#### 1.2 ホメオログ遺伝子の発現量定量化

本研究では、遺伝子発現量の定量化に RNA-seq を用いる。RNA-seq は mRNA の配列断片を大量に読み (配列をリードと呼ぶ)、既知の遺伝子配列やゲノム配列と対応したリードの本数を数えることで遺伝子の発現量を計測する。RNA-seq をホメオログ遺伝子発現量の定量化に用いることを考えよう。 $c_a, c_b$  由来のリードは、a, b の配列両方に対応するため、由来する遺伝子を特定できず、発現量の定量化を誤る可能性がある。本手法では、ホメオログ遺伝子間の変異を利用して各リードがいずれの親由来であるか分離する。

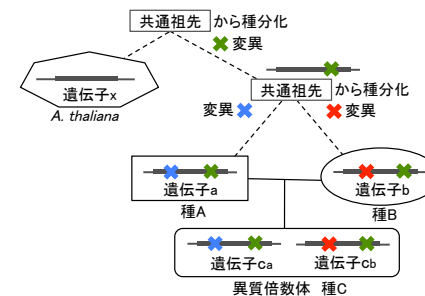


図 1 系統樹、遺伝子の持つ変異の関係

Fig. 1 Relationships between Mutations and Phylogenetic Tree

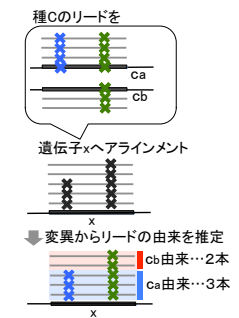


図 2 手法の概要

Fig. 2 Outline of Proposed Method

## 2. 手 法

RNA-seq では、リードをアラインメントするゲノムが必要だが、異質倍数体と親種のゲノムは未知のため、近縁で、ゲノムや遺伝子配列が既知の *A. thaliana* を利用する。

### 2.1 変異の種類と定義

a, b を *A. thaliana* の遺伝子  $x$  のオソログとする。系統樹を考えると、 $x$  と  $c_a, c_b$  の配列間の差異は、親種特有の変異 (図 1 の  $c_a, c_b$  上の  $x$  印, 赤青), 祖先種から親種へ進化する過程で起きた変異 (緑)(以下, 親由来の変異, 祖先由来の変異) の 2 種類に分けられる。 $c_a, c_b$  由来のリードを  $x$  にアラインメントした模式図が図 2 である。

ある一つの遺伝子上に変異が  $k$  箇所あるとき、変異を有する位置の集合を  $M = \{M_1, M_2, \dots, M_k\}$  とする。 $M_i$  の位置を含むリードの集合を  $R_i$ , 要素数を  $|R_i|$  と表す。変異の有無で分けたリードの集合を  $R_{1i}, R_{2i}$  とする。この時、 $|R_{1i}| \leq |R_{2i}|$  とする。

### 2.2 由来別の遺伝子発現量の算出

図 2 のように、アラインメントした時点では変異の種類は分からず、変異が、親由来か祖先由来かを見分ける方法が必要になるが、後述する。計算により判定した親由来の変異に注目すると、リードは変異の有無から 2 本と 3 本に分けられ、リードごとにいずれの親由来か (青, 赤) が分かる。更にリードの本数から、リードの由来別のホメオログ遺伝子の発現量が求められる。ここでは特定の変異に着目して親種の比率を求めたが 1 点のみの推定ではシーケンシングエラーによる影響も多く精度が低い。本研究では、リード毎にいずれの親由来であるかを決定し、その本数を遺伝子全体にわたって数える事で計算の精度を向上させる。

### 2.3 変異の種類を見分ける方法

変異  $M_i$  に着目した時、この変異は親由来か祖先由来のどちらかである。RNA-seq はランダムに配列を選択する手法であり、 $M_i$  のリードの本数が  $R_{1i}, R_{2i}$  である確率  $P(M_i)$  を、ポアソン分布を用いてモデル化する。 $\lambda_p$  を変異が親由来である時の、変異のあるリードの本数の期待値、同様に  $\lambda_a$  を祖先由来の変異である時の期待値とすると、 $M_i$  が親由来である確率は、 $P(M_i) = \lambda_p^{|R_{1i}|} e^{-\lambda_p} / |R_{1i}|!$ , 祖先由来である確率は、 $P(M_i) = \lambda_a^{|R_{1i}|} e^{-\lambda_a} / |R_{1i}|!$  と表せる。 $t$  を祖先由来の時の変異率とすると、 $\lambda_a = t \times |R_i|$  となる。 $t$  はあらかじめ固定した値である。 $\lambda_p$  は遺伝子ごとに異なる値であり、データから求める。

それぞれの変異が独立な事象とすると変異  $M$  の尤度  $P(M)$  は  $P(M_1) \times P(M_2) \times \dots \times P(M_k)$  で表せるので、この尤度が最大になるように、すべての変異を親由来か祖先由来かに分ければ良い。この時、親由来と考えられる変異セットを  $M_P (\subseteq M)$  とする。

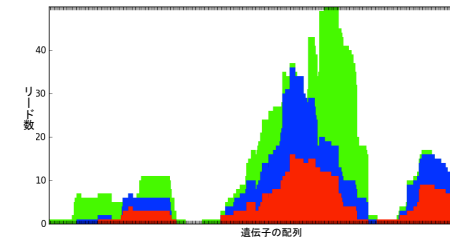


図 3 親種を区別したリード本数の分布

Fig. 3 Distribution of Number of Reads Classified by Originated Parental Species

$M_P$  を現実的な時間で求めるため、以下の工夫を行った。

- (1)  $M_i \in M_P$  に対して  $\lambda_p$  は  $M_P$  に属する変異の平均に近いと考えられるので、 $I = \{i \mid M_i \in M_P\}$  とすると、 $\lambda_p = |R_i| \times (\sum_{i \in I} |R_{1i}| / \sum_{i \in I} |R_i|)$  とする。
- (2)  $|R_{1i}| / |R_{2i}|$  が大きいもの程、 $\lambda_p$  を大きくするので、この比率が 1 に近いものから順に  $M_P$  に入れて、尤度が最大となる組み合わせを計算した。

## 3. 実データを用いた実験の結果

異質倍数体の *C. flexuosa* に対し Life Technologies 社の SOLiD5 を用いて RNA-seq を行ったリード約 6 千万本に対し、本手法を適用した。アラインメント先に *A. thaliana* (TAIR10)<sup>1)</sup> のゲノムを、アラインメントソフトには SHRiMP<sup>2)</sup> を用いた。図 3 は、一つの遺伝子に対し本手法を適用した結果を示している。横軸に遺伝子上の位置、縦軸に塩基ごとのリードの本数を表す。色分けは、リードごとに判定した親種に対応し、赤と青が片親、緑は由来が不明のリードである。この遺伝子の両親種ごとの遺伝子発現量は、それぞれ 5.93, 6.56 だった (RPKM<sup>3)</sup> 換算)。グラフからは、この遺伝子の場合、親種ごとの発現量を求めるのに使うことの出来たリードがほぼ遺伝子の全体を網羅し、発現の比が遺伝子上を一貫して一定であったことが分かる。

## 参 考 文 献

- 1) Lamesch, P. et al. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Research*, **40** (D1): D1202-D1210., 2012.
- 2) Rumble, S.M. et al. SHRiMP: accurate mapping of short color-space reads. *PLOS Comput. Biol.* **5**, e1000386, 2009.
- 3) Mortazavi, A. et al. Mapping and quantifying mammalian transcriptomes by rna-seq. *Nature Methods*, **5**(7):621-628, 2008.