



モンテカルロ法による分岐図 作成のための最適化法

東京農工大学大学院

連合農学研究科

松田 朋子

次世代シーケンサーによる遺伝子解析 大量のデータをどう扱うか？

サンガーシーケンサー



次世代シーケンサー



リード長 約 1,000 bp
1 ランあたりの時間 2 時間
1 ランあたりのデータ量
約 0.016 Mb

リード長 100 bp
1 ランあたりの時間 約4日間
1 ランあたりのデータ量
約 95,000 Mb

解読量が大幅に増加



データ解析が困難



目的

どの遺伝子を使うのか？

- 従来，分岐図は，数個の遺伝子（ミトコンドリアDNAの一部の配列など）をもとに作成していた
- 近年のゲノム解析技術の進歩にともない，ゲノム全体にわたる**大量の遺伝子情報を容易に得られる**ようになった
- 大量の遺伝子情報のうち，**どの遺伝子**を使って分岐図を作成するのかを検討する必要性が出てきた
- 今回は，大量のデータから無作為にサンプルを抽出する方法（モンテカルロ法）を用いて，**無数の遺伝子の組合せ**を作る方法を確立することを目的とした



方法

1. 組合せ生成

- 447種類の遺伝子の組合せ数を計算して、組合せのパターンを生成
- 447種類の遺伝子のうち1~4遺伝子を除く組合せの数は“1656133808”

2. 乱数と数列の作成

- “1656133808”から一様乱数を発生させる
- 乱数に対応する組合せ数列を作成

3. 分岐図作成

- RAxMLによる分岐図作成

4. 分岐パターン抽出

- 系統樹を推定し、分岐パターンを抽出



結果①：組合せ生成

447個の遺伝子から任意の個数の遺伝子を選択する組合せの数を計算

- `> CHOOSE(447,1)`
[1] 447
- `> CHOOSE(447,1)+CHOOSE(447,2)`
[1] 100128
- `>CHOOSE(447,1)+CHOOSE(447,2)+CHOOSE(447,3)`
[1] 14886143
- `>CHOOSE(447,1)+CHOOSE(447,2)+CHOOSE(447,3)+
CHOOSE(447,4)`
[1] 1656133808 <- 本研究で使用した乱数の数. 選択した
1~4個の遺伝子を解析から除く



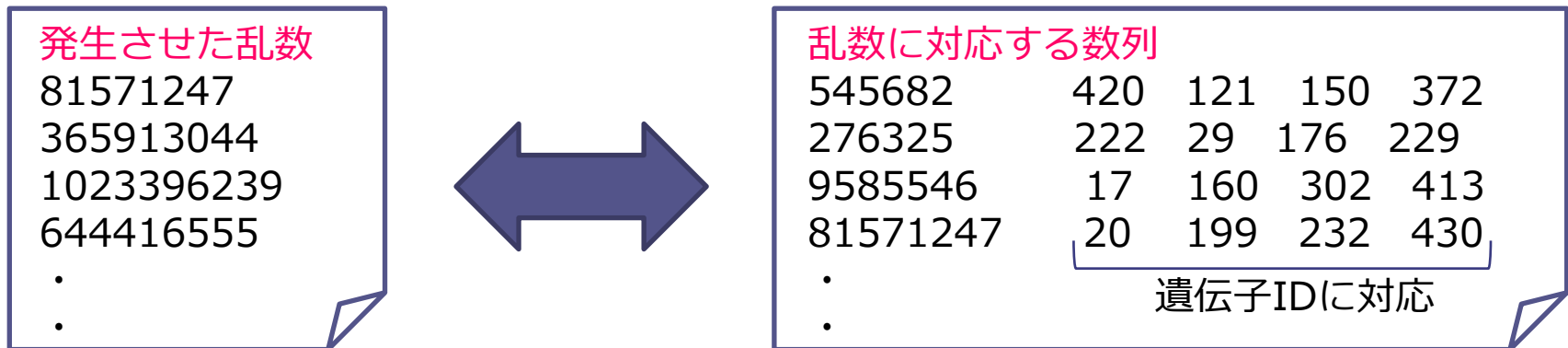
結果②：乱数の発生と対応する組合せ数列の選択

1. 一様乱数

```
RUN1 <- ROUND(RUNIF(100000, MIN = 1, MAX = 1656133808))
```

2. 447個の遺伝子から4個までを選択する1656133808個の組合せ数列を 発生

3. 乱数と数列を対応させる



4. 447個のうち乱数に対応する遺伝子を除き，残りのすべての遺伝子をつかって分岐図を作成



結果③ : RAxMLによる分岐図作成

- 分岐図作成ソフトRAxMLを実行
- RAxML の出力結果から数値や記号などを取り去り, 分岐図パターンをシンボリックデータにする

RAxML の出力結果

```
((x:0.00531404813722460411,b:0.00582680929801783314):0.02597126625043840939,(m:0.03645324955994154459,(i:0.02130122160079299734,g:0.03144790765745542060):0.00424151281867054565,u:0.02948686769656536782):0.02040360046940021752):0.01521774994899611003):0.00840806219060770237,(((z:0.34344767189954156228,(t:0.14262613954560879326,l:0.09611024306570406517):0.12482727188754781655):0.17738679389459199864,((q:0.09220903636333667441,c:0.07849740402964605623):0.02162213710044687265,(((e:0.02209622018870339294,k:0.02893283119099681472):0.00897154535047478552,p:0.08526167426794276083):0.01393193949473354662,(f:0.01318356630444799359,(n:0.01952490523202302791,&:0.01693127308779425119):0.00813604928815400003):0.07817988855466992404):0.02263432315084732208,(r:0.04590445767519640841,h:0.04637315388999797144):0.03865318679664145329):0.00789191373814705256,(v:0.03284611917996409919,d:0.02250210568318532570):0.17923032798411692168):0.01001049487098192720):0.10325734189742048763,(y:0.20086249354886381857,(j:0.18321752975378832740,#:0.20967985390326032702):0.12992858329809614526):0.01880523845322217696):0.03857414591748902638):0.06620101031901222399,(o:0.03854648520093560682,s:0.03552704927452698946):0.12706855170728659221):0.05630830036902149949,w:0.13933629626394547496):0.07157780409461377003,a:0.03782021720159066402):0.0;
```

⇒簡素化した分岐図パターン

```
((x,b),(m,((i,g),u))),(((z,(t,l)),((q,c),(((e,k),p),  
,(f,(n,&))),),r,h)),(v,d)),(y,(j,#))),),o,s)),w),a)
```



結果④：分岐図パターンの頻度

185/200 試行 で得られた樹形

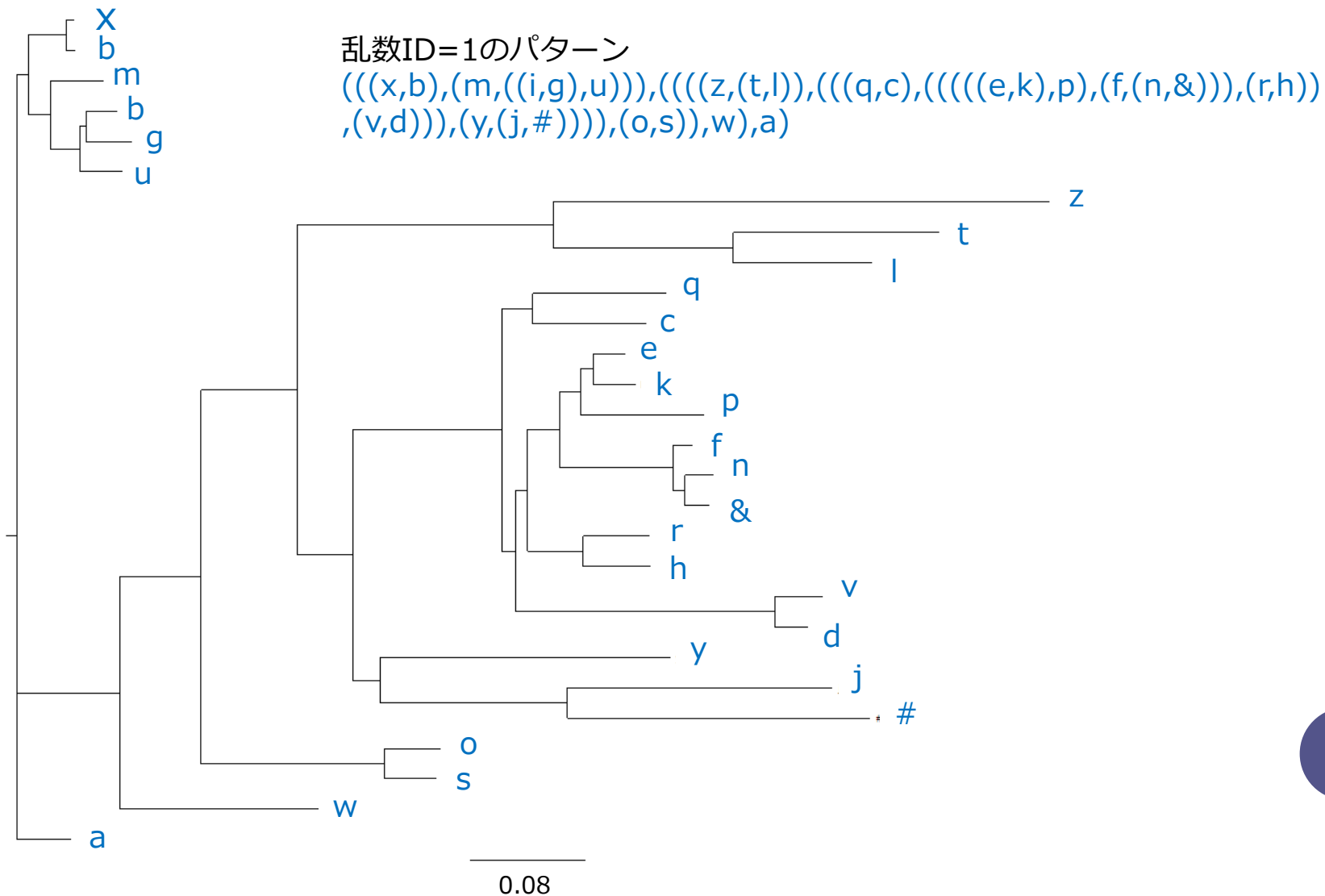
`((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)`

15/200 試行 で得られた樹形

`((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),((v,d),(((e,k,p),(f,(n,&))),(r,h))),(y,(j,#))),(o,s)),w),a)`

Random Number	ID	Gene1	Gene2	Gene3	Gene4	Tree Pattern
all						<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
1	5	121	150	372		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
2	27	29	176	229		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
3	95	160	302	413		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
4	51	199	232	430		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
5	56	126	142	201		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
6	135	342	359	446		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
7	89	303	405	437		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
8	54	98	149	213		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
9	39	234	259	417		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
10	84	108	135	403		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
11	107	277	308	420		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
12	132	379	399	432		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
13	20	223	367	414		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
14	23	96	240	421		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
15	144	147	208	414		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
16	77	267	273	309		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),((v,d),(((e,k,p),(f,(n,&))),(r,h))),(y,(j,#))),(o,s)),w),a)</code>
17	20	243	277	295		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
18	6	219	248	317		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
19	17	73	271	274		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>
20	56	181	237	313		<code>((x,b),(m,((i,g,u))),(((z,(t,l)),((q,c),(((e,k,p),(f,(n,&))),(r,h)),(v,d))),(y,(j,#))),(o,s)),w),a)</code>

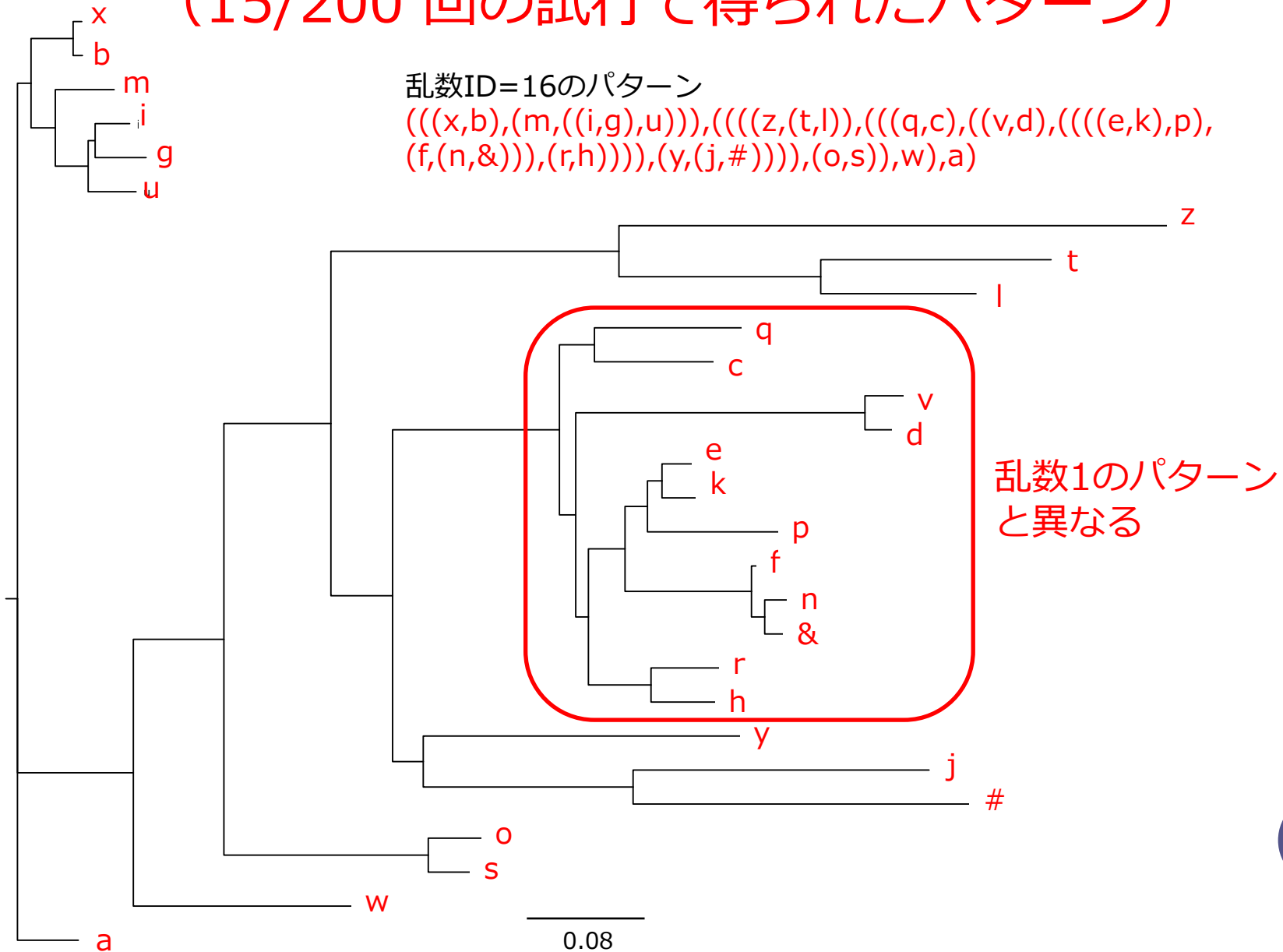
結果⑤最適化された分岐図パターン (185/200 回の試行で得られたパターン)



結果⑥最適化された分岐図とは異なるパターン (15/200 回の試行で得られたパターン)

乱数ID=16のパターン

`((x,b),(m,((i,g),u))),(((z,(t,l)),((q,c),(v,d),(((e,k),p),
(f,(n,&))),r,h))))),(y,(j,#))),o,s)),w),a)`



乱数1のパターン
と異なる



まとめ

1. 大量のデータから無作為にサンプルを抽出する方法（モンテカルロ法）を用いて、**無数の遺伝子の組合せ**を作る方法を確立した
2. 200回の試行のうち、185回の試行で完全に一致するパターンが得られた。一方、残り15回の試行では、異なるパターンが得られた。
3. 今後は、より多くの試行を実施して、異なるパターン（15/200回の試行）が得られた時に除かれた遺伝子（=185/200回の試行で得られたパターンを得るために必要な遺伝子であると考えられる）をリストアップする。リストアップした遺伝子の中には、**分岐図作成に強く関与**している遺伝子が含まれる可能性が高いと考えられる。

