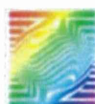


# Garuda Alliance (snapshot May 2013)



TTG4 菅野班 H26年度 総括 2012-02-13 version

63

## まとめ

- 新型反復暴露解析:
  - 過渡反応と基線反応の基本的な関連性を見いだした。毒性学的に新規性が高くエピジェネティクス等の機序の関与が示唆される。上流にEif2、RICTORなどが想定された。
  - 反復毒性の分子毒性学的理解の促進、及び、単回暴露データからの反復毒性予測に重要。
- 胎児発生過程におけるマスター遺伝子を基軸とした遺伝子発現ネットワークの網羅的解析:
  - 微分解析手法により、網羅的・効率的に発生マスター遺伝子の解析に成功した。
  - 未報告の上流制御系が抽出される可能性が見いだされた。
- Percellome 3次元データ等の為の専用解析ソフトウェアの開発研究
  - 異種動物データの統合技術
  - 絶対量化されていない非Percellomeデータの絶対量推定技術
  - RSort改良
  - WebAPI改良
  - その他
- システムトキシコロジー解析基盤の研究開発:
  - 状態制御遺伝子群推定アルゴリズムAGCT
  - オルソログ制御領域分析ソフトウェアSHOE
  - 遺伝子制御関係推定アルゴリズム
  - Web公開/Garuda Platform実装/AOP適合/Networkによるスクリーニングパイプライン
  - 次世代Computational Tox Screening Systemの創出

TTG4 菅野班 H26年度 総括 2012-02-13 version

64

## 委託研究報告書 (STEP13)

次世代シーケンサを利用した遺伝子発現解析の高度化～スプライシングバリエントの識別アルゴリズムの改良

Copyright(C)2014-2015 NTT DATA Corporation

# 1.平成26年度研究テーマ

次世代シーケンサを利用した遺伝子発現解析の高度化

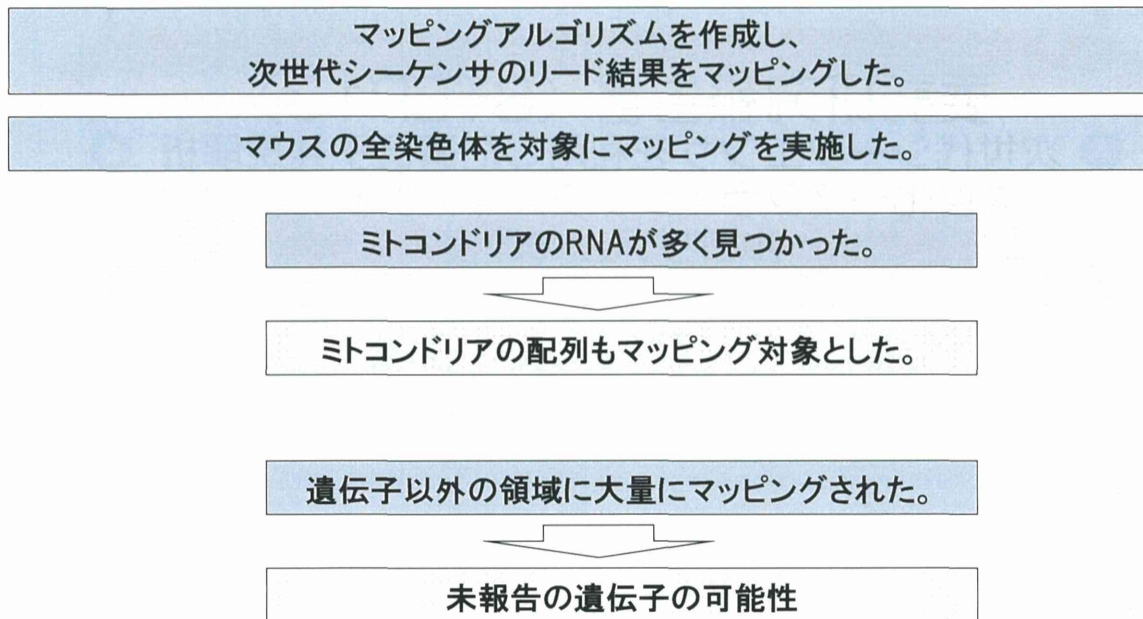
スプライシングバリエントの識別アルゴリズムの改良

先行研究より、次世代シーケンサを利用した遺伝子発現解析の高度化は、マッピングアルゴリズムの性能向上が欠かせないことが明らかになってきた。特に、未報告のエクソン構造・スプライシングパターンを持つメッセンジャーRNAの検出・定量には、高度なマッピングアルゴリズムの開発が必須である。

本研究では、マッピング精度を高め、スプライシングバリエントの推定・識別を可能とするための技術開発を行う。

Copyright(C)2014-2015 NTT DATA Corporation

## 2.1.昨年度の成果



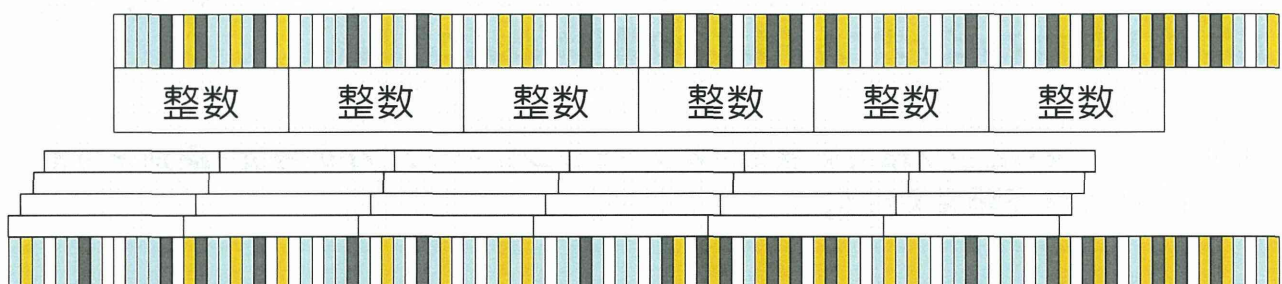
Copyright(C)2014-2015 NTT DATA Corporation

## 2.2.Teradataを用いたマッピングの基本アイデア

読込配列と参照配列を、15塩基ずつに分割し、それぞれ、30ビット(4バイト)整数で表現する。

Teradata RDBMSは、ハッシュインデックスを用いて、完全一致の検索を高速に実行可能である。

### 計測配列



### 参照配列

1塩基ずらしたパターンを生成する。

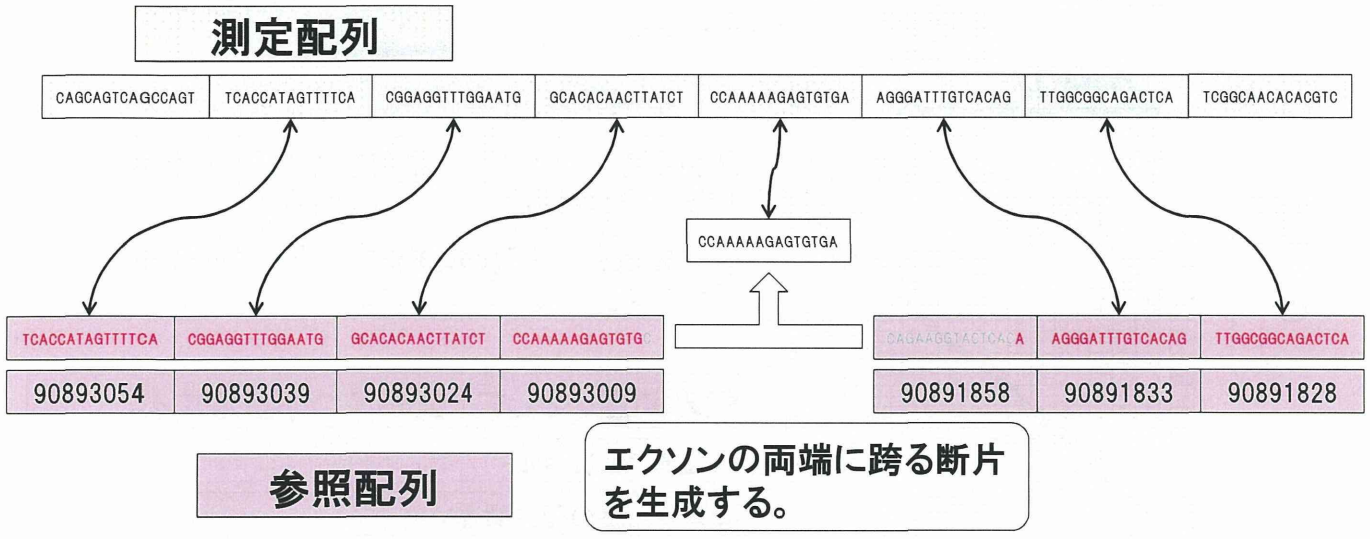
スプライシング前と後の2種類を生成する。

Copyright(C)2014-2015 NTT DATA Corporation

## 2.3. スプライシング後のRNAをマッピングさせる方法

### エクソンに跨る配列のマッピングマッピング方法

15塩基の断片の情報と、どの断片が隣接しているかの情報のみを蓄積している。エクソンに跨る断片を疑似的に生成し、隣接情報を生成する。

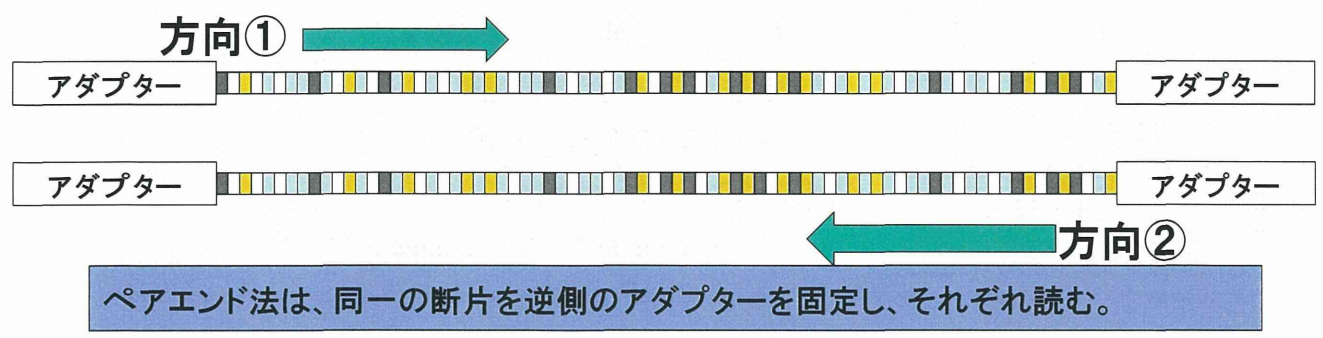


Copyright(C)2014-2015 NTT DATA Corporation

## 2.4. ペアエンド処理

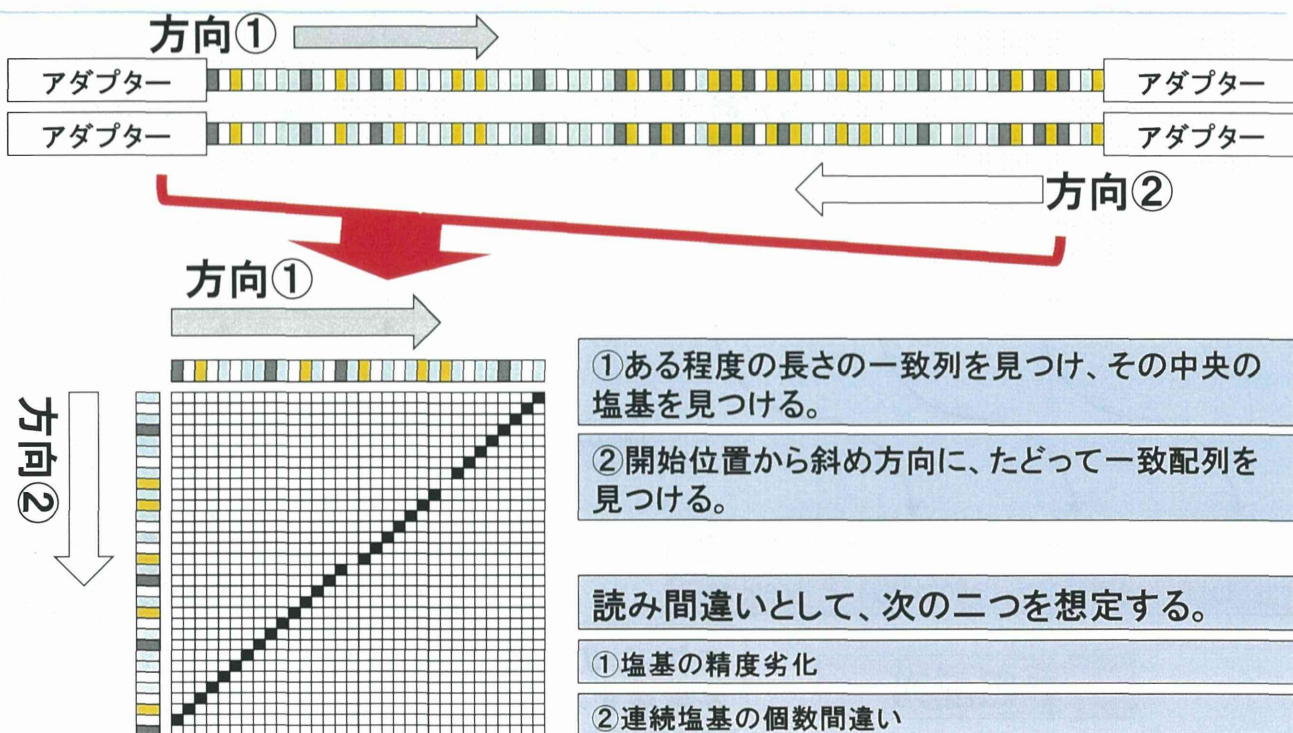
ペアエンド法の両側から読むことで、信頼性の高い配列を推定する。

ペア両方のファイルを読み込んで、処理しながらTeradataにロードしていく。



Copyright(C)2014-2015 NTT DATA Corporation

## 2.5.ペアエンド処理：アルゴリズム概要

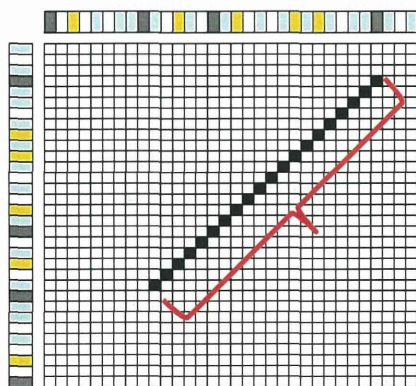


Copyright(C)2014-2015 NTT DATA Corporation

## 2.6.ペアエンド処理：アルゴリズム①

アダプター検索で場合分けし、一致配列を探す開始位置を見つける。

制約条件なしに一致箇所を探そうとすると、非常に時間がかかる。  
ある程度の精度で読み取れているのならば、ある程度の長さで一致しているはずである。  
そこを起点として調べる。



両方でアダプターが見つかった。

長さ3未満の違い

長さを変えながら、一致配列を探す。

長さ3以上の違い

諦める。

片方だけでアダプターが見つかった。

長さを変えながら、一致配列を探す。

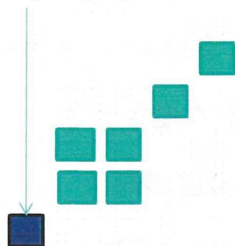
両方ともアダプターが見つからなかった。

30~400の範囲で一致の中心箇所を探す。

Copyright(C)2014-2015 NTT DATA Corporation

## 2.6.ペアエンド処理：アルゴリズム②

(a,b)までマッチしているとする



(a+1,b+1)をチェックする。

一致すれば、他はチェックせず、進む。

一致しなければ、3か所チェックし、分岐する。

(a+2,b+2)をチェックする。

一致すれば、コスト+1で、進む。

一致しなければ、対角線上でチェックしつづける。

(a+2,b+1)をチェックする。

一致すれば、コスト+3.4で、進む。

一致しなければ、この分岐を諦める。

(a+1,b+2)をチェックする。

一致すれば、コスト+3.4で、進む。

一致しなければ、この分岐を諦める。

3分岐のコストが一番安いものを採用する。

対角できれいに並んでいる場合には、分岐せずに、コスト0で進める。

不一致と、塩基ずれの同時発生は諦める。

Copyright(C)2014-2015 NTT DATA Corporation

## 2.7.ペアエンド処理の効果

アダプターの読み間違いによる長さ不定による棄却の減少

塩基の読み間違いによる非マッチの減少

➡ 読取品質の向上

振動などによる品質低下に対する抑止

振動などで特定位置での読取品質が低下する可能性がある。

読み始めから、何塩基目かで集中する。

逆方向からは、長さは集中しない。

品質の良い方を採用することで、全体の品質を向上させる。

### 3.1.今年度の方針

#### 1. マッピング精度の向上

非マッチリード配列の検証

- エクソン情報の追加
- 参照配列の補正 など

複数の参照配列を使用できるようにする。

- 新たな知見で、配列を修正する。
- マッピング精度の向上
- 対象配列を絞り込むことで実施する。
- マッピングの高速化

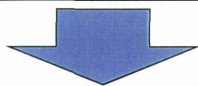
#### 2. マッピングの高速化

対象配列を絞り込むことで実施

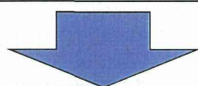
### 3.2.高精度化のためのマッピングエラー配列の特定

#### 特定方法①

大量に存在するはずの遺伝子      Alb、ApoEなど



エクソン中の塩基配列      30塩基程度



一致したリード配列で、マッピングできなかった配列を見つけ出す。

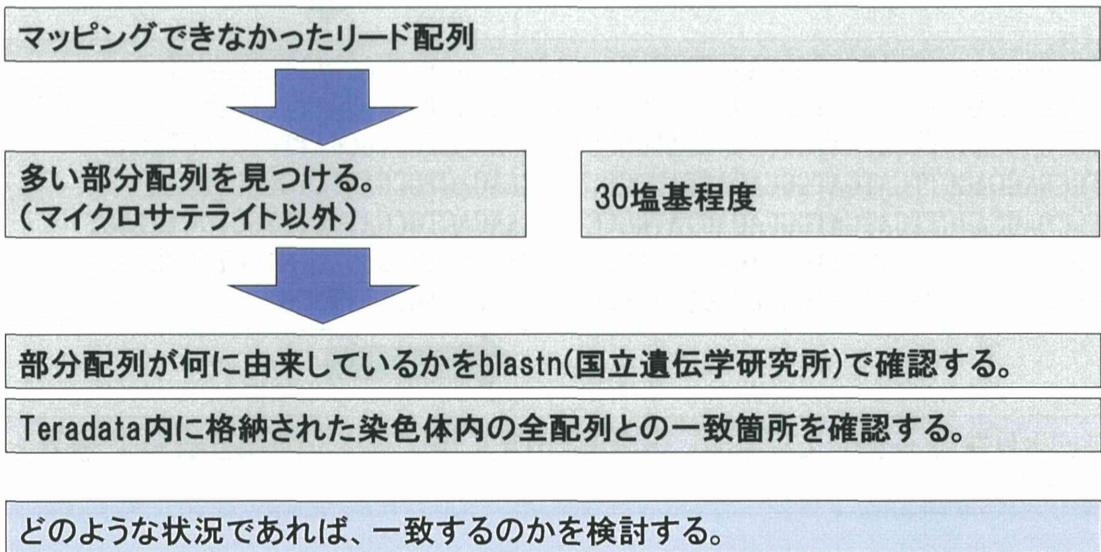
どのような状況であれば、一致するのかを検討する。





### 3.4.高精度化のためのマッピングエラー配列の特定

#### 特定方法②



Copyright(C)2014-2015 NTT DATA Corporation

### 4.5.高精度化のためのマッピングエラー配列の特定

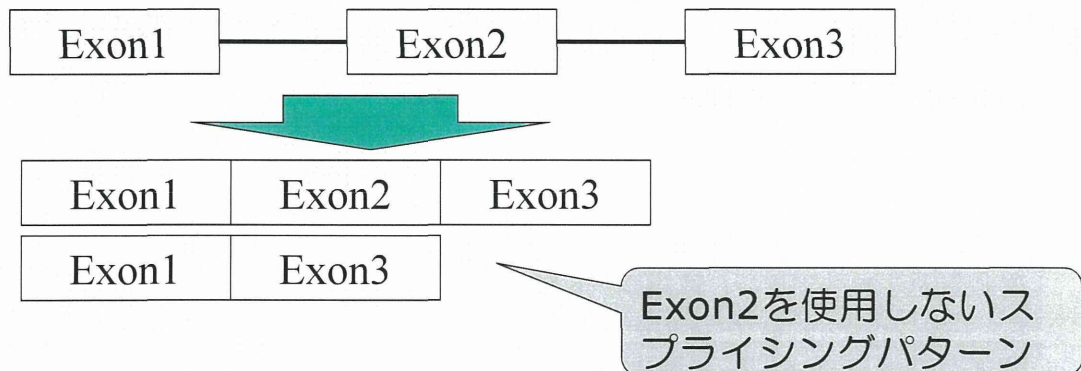
#### マッチしなかったリード配列の部分配列が多かった上位30

No	配列	個数	blastn description
1	CTTAGGTATAGTAAATGATTGAATCCATCATACGTCACAGAATA	14340	Major Urinary Protein MUPS MUP11 chr4
2	AGAATATTCACCCAGCCTTTTCTGTTTTGTTCAGCAACCATAGATAA	12217	Major Urinary Protein MUPS MUP11 chr4
3	AGGAAGGGATGATGGTGGAGCTCGGTGAGAAGTCTCCACTCAAC	11338	Major Urinary Protein MUPS MUP11 chr4
4	GTCCTCACTCAACACTGGAGGCTCAGGCCATTCTTCATTCTCGGG	9811	Major Urinary Protein MUPS MUP11 chr4
5	GCTGGAGTCTGGTGAAGAAGTCTCCACTCAACACTGGAGGCTCAG	9531	Major Urinary Protein MUPS MUP11 chr4
6	TTGGCATTGGATAGGTCAATGATATTTTCTCTAAGGATTCATGC	8315	Major Urinary Protein MUPS MUP11 chr4
7	CTFGAACTCTCTCTCAAAAGTCTTTTCAACTTTCCCTTACGGTA	7885	45S pre rRNA gene / 28S ribosomal RNA
8	GCCCTCTTGAACCTCTCTTCAAAGTCTTTTCAACTTTCCCTTA	7877	45S pre rRNA gene / 28S ribosomal RNA
9	CAGAATATTCACCCAGCCTTTTCTGTTTTGTTCAGCAACCATAGATA	7654	Major Urinary Protein MUPS MUP11 chr4
10	CTGGGATGCTGTATGGATAGGAAGGGATGATGGTGGAGTCTCGG	6890	Major Urinary Protein MUPS MUP11 chr4
11	CTCTTGAACCTCTCTCTCAAAGTCTTTTCAACTTTCCCTTACGG	6511	45S pre rRNA gene / 28S ribosomal RNA
12	GAATATTCACCCAGCCTTTTCTGTTTTGTTCAGCAACCATAGATAA	6406	Major Urinary Protein MUPS MUP11 chr4
13	GCATTGGATAGGTCAATGATATTTTCTCTAAGGATTCATGCTCC	6259	Major Urinary Protein MUPS MUP11 chr4
14	CCATCAGCTGGAAGGTTTCCCATCCTTTTCGTTAATGAGATGAG	5808	Major Urinary Protein MUPS MUP11 chr4
15	GTCACAGAATATTCACCCAGCCTTTTCTGTTTTGTTCAGCAACCAT	5803	Major Urinary Protein MUPS MUP11 chr4
16	GGAGGGGGCCGGCCGCCACCCACCCAGCCCGCCGGGAGGCGG	5784	45S pre rRNA gene / 28S ribosomal RNA
17	AATATTCACCCAGCCTTTTCTGTTTTGTTCAGCAACCATAGATAA	5585	Major Urinary Protein MUPS MUP11 chr4
18	CTCTCTTCAAAGTCTTTTCAAAGTCTTTTCAACTTTCCCTTACGG	5449	45S pre rRNA gene / 28S ribosomal RNA
19	CCGGAGGCGGAGCGGGGGAGAGGGAGAGCGGCGGCGAGGGTATC	5447	45S pre rRNA gene / 28S ribosomal RNA
20	CAAACTTTCTTCTGATGCTGAACTCAAATCTGGTTCTCGGCCAT	5275	Major Urinary Protein MUPS MUP11 chr4
21	CGGGAGCGGAGCGGGGGAGAGGAGCGCGGCGAGCGGGTATCT	5202	45S pre rRNA gene / 28S ribosomal RNA
22	GTGAGAAGTCTCACTCAACACTGGAGGCTCAGGCCATTCTTCAT	5084	Major Urinary Protein MUPS MUP11 chr4
23	GCCATTATCTCTATCTTTCTCTTTTGTTCAGAGGCCAGGATAAT	5070	Major Urinary Protein MUPS MUP11 chr4
24	TGCCATTATCTCTATCTTTCTCTTTTGTTCAGAGGCCAGGATAA	5042	Major Urinary Protein MUPS MUP11 chr4
25	TCCGACCTGGGCGGGTTCACTCCTCTTAGGCAACTCGGTGGTC	4993	Mus musculus clone 2S-385 MMuLV retroviral integration site genomic sequence.
26	CACTGGAGGCTCAGGCCATTCTTCATTCTCGGCCCTGGAGGCGC	4784	Major Urinary Protein MUPS MUP11 chr4
27	CCCTCTTGAACCTCTCTCAAAGTCTTTTCAACTTTCCCTTAC	4728	45S pre rRNA gene / 28S ribosomal RNA
28	GTCTTAAGGCAGCTCAGGGAGGACAGAACTCCCGTGGAGCAG	4726	45S pre rRNA gene / 28S ribosomal RNA
29	GAAACGGGAGCGGGAAGATCCGCCGGGACACCGGCACGGCCG	4720	45S pre rRNA gene / 28S ribosomal RNA
30	GGATGATGGTGGAGTCTGGTGAAGAAGTCTCACTCAACACTGGA	4654	Major Urinary Protein MUPS MUP11 chr4

MUPとRibosomalRNAが多い、配列、エクソン情報の誤りの可能性がある。

## 4.6.昨年度のアルゴリズムの未達成事項

### スプライシングバリエーションの想定



いくつかの遺伝子では、エクソンパターンが多く存在しており、このようなエクソンを使用していないパターンも網羅しているケースも存在した。

⇒ エクソンを飛ばした、跨りのパターンを追加する。

## 5.参照用配列

参照配列として、以下の配列を用いる。

- ・ ダウンロードサイト
  - Cufflinksサイト
  - <http://cufflinks.cbcb.umd.edu/igenomes.html>
  - Mus musculus/UCSC/mm10      Jun 14 11:29
- ・ 対象
  - Mm10内の情報
  - 総配列
    - ・ Chromosome
  - 遺伝子アドレス情報
    - ・ Ref\_text