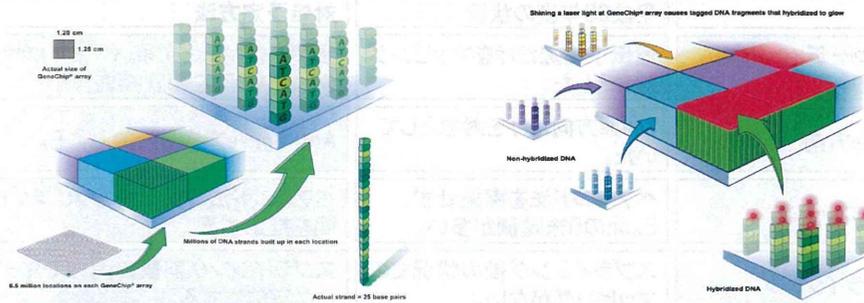


マイクロアレイ上にRNAの各領域に相補鎖を作成し、対象サンプルを蛍光で標識する。その領域でどの程度の蛍光輝度があるかで、結合したRNA量(発現量)を測定する。

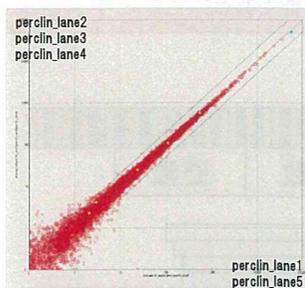


マイクロアレイの問題点

- ・未知の配列を検出できない。
- ・飽和現象を起こしやすい。
- ・類似配列への誤結合(クロスハイブリダイゼーション)を起こすことがある。

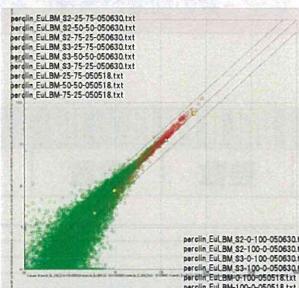
試験管内で混合した25%:75%、50%:50%、75%:25%を数値的に平均をとった値を縦軸に、Liver100%とBrain100%を数値的に平均をとった値を横軸にして散布図を作成した。

シーケンサ



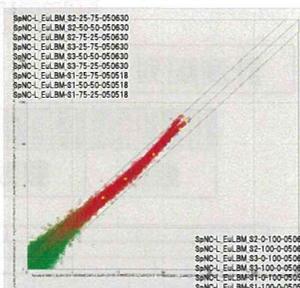
高発現から低発現まで
対角線上にある

マイクロアレイMLANG



高発現は対角線上に
乗っている

マイクロアレイMAS5



高発現で分岐する状況がある

マイクロアレイでは飽和現象があり、MAS5では、Liver側とBrain側で絶対量が異なるため、飽和パターンが異なり、分岐したと思われる。MLANG(特許第517712号)では、飽和をキャンセルできたが限界がある。

平成23年度委託研究で明らかになった課題と対処予定方法

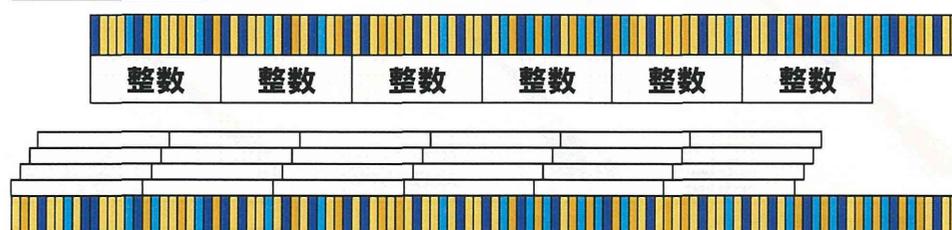
項目	平成23年度の状態	対処予定方法
ノンコーディングRNAへの対応	遺伝子周辺だけをマッピング対象とした。	全DNAを対象とする。(マイクロサテライトの排除を適宜実施)
Sense方向	Sense方向だけを対象としていた。	Anti-Sense方向も対象とする。
ペアエンド法	ペアエンド法を実施せず、Exonの5末端側が多い。	ペアエンド法の両端をマッピングし、間を推定する。
スプライシング	スプライシング後の状況でのマッピングがない。	スプライシング前後の両方でマッピング可能にする。
複数箇所へのマッチ	複数箇所にマッチした場合は平均按分。	周囲の状況を考慮した按分を実施する。

3. 1. Teradataによるマッピングの基本アイデア①

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は、完全一致の検索が高速に実施可能である。

計測配列



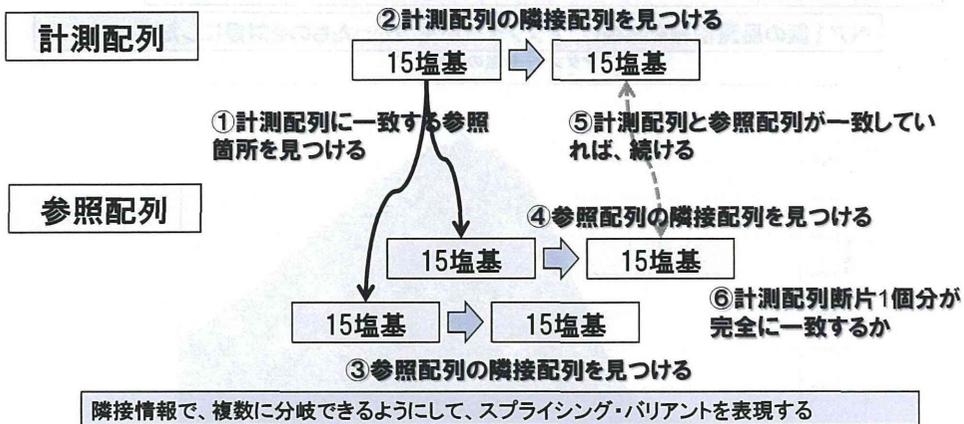
参照配列

1塩基ずらしたパターンを生成する

全塩基数の情報を生成する

スプライシング前と後の2種類を生成する

①→⑥の順番で処理を実施することにより、計測断片全長で、参照配列上と一致するかを判定可能である



課題点

- 参照配列上に同じパターンが多数存在すると、処理しきれなくなる
- 完全一致以外では、工夫が必要

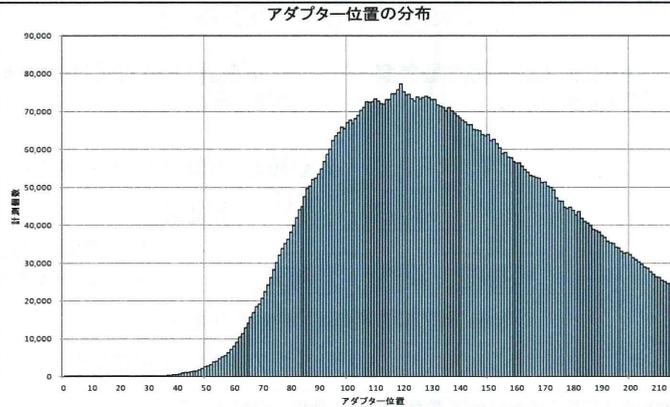
参照配列上に、未定塩基(N)である箇所が何箇所が存在する。

連続個数	存在数
1	47
2	8
3	2
4~7	0

3連続以下の未定塩基は対処すべきである。8個以上連続の未定塩基列へのマッピングは困難であり、今回は対象外とする。

測定断片の長さ(アダプター位置)の分布を確認した。

ペア側の品質情報を無視し、アダプターが見つかったものを対象にした。



長さ120前後でピークとなる分布形状であった。

マッピングアルゴリズムの基本部分のみを作成し、課題の洗い出しを行った。

対象

ヘアエンドの一方だけを使用

ランダムサンプル1/10000

対象配列数

1009

計測配列中に未定塩基を含まない
アダプターが見つかった

対象配列数

798

マッピングできた配列

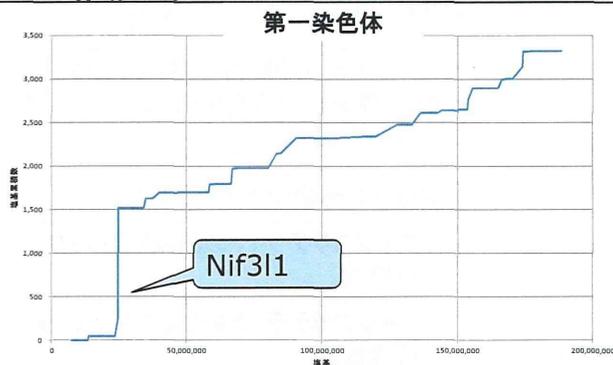
219

マッピングできたペア(累計)

2057

4. 2. 第一染色体におけるマッピング結果

マッピングした配列を、塩基アドレスと、マッピング累積塩基数の積み上げ数でグラフを作成した。



マッピング結果のペアで、測定配列と参照ゲノムの配列は一致しており、基本的なアルゴリズムが動作している事が確認できた。
Nif31 近傍 (遺伝子領域外) でマッピングされるリードが集中して表れた。

4. 3. マイクロサテライト

短い配列の繰り返しは、ゲノム上に大量に存在し、マッピング処理の処理性能を著しく低下させる。マイクロサテライトの状況を確認した。

	ゲノム上での存在数	配列パターン	基本パターン番号	マイクロサテライトパターン	開始位置
1	3270514	ACACACACACACA		6 CA	1
2	3269104	TGTGTGTGTGTGT		11 TG	0
3	3184366	CACACACACACAC		6 CA	0
4	3182110	GTGTGTGTGTGTG		11 TG	1
5	1792080	TCTCTCTCTCTCT		15 TC	0
6	1782146	AGAGAGAGAGAGA		5 GA	1
7	1761334	CTCTCTCTCTCTC		15 TC	1
8	1751890	GAGAGAGAGAGAG		5 GA	0
9	1083204	TTTTTTTTTTTTT		3 T	0
10	1073202	AAAAAAAAAAAAA		0 A	0
11	958494	ATATATATATATA		7 TA	1
12	957864	TATATATATATAT		7 TA	0
13	391122	AAAGAAAGAAAGAA		85 GAAA	1
14	387138	TTTCTTTCTTTCTT		275 TTTC	0
15	382322	AGAAAGAAAGAAAG		85 GAAA	3
16	379226	AAGAAGAAAGAAAG		85 GAAA	2
17	379226	GAAAGAAAGAAAGAA		85 GAAA	0
18	377732	TCTTTCTTTCTTCT		275 TTTC	2
19	374820	TTCTTTCTTTCTTC		275 TTTC	1
20	374820	CTTTCTTTCTTCTT		275 TTTC	3

➡ 測定配列全体が、マイクロサテライトであるならば、その測定断片は、対象外とする。

4. 4. マイクロサテライト以外の非常に多い配列

No.	件数	配列パターン
1	122743	GTGTGTGTGTGTGTA
2	122743	ATGTGTGTGTGTGTG
3	122143	TACACACACACACAC
4	122143	CACACACACACACAT
5	93657	TGTGTGTGTGTGTAT
6	93657	TATGTGTGTGTGTGT
7	92357	ATACACACACACACA
8	92357	ACACACACACACATA
9	80888	CTTTTTTTTTTTTTTT
10	80888	TTTTTTTTTTTTTTTC
11	80257	GAAAAAAAAAAAAAAA
12	80257	AAAAAAAAAAAAAAG
13	79363	GTTTTTTTTTTTTTTT
14	79363	TTTTTTTTTTTTTTTG
15	79329	AAAAAAAAAAAAAAAC
16	79329	CAAAAAAAAAAAAAAA
17	73675	AACAACAACAACAACA
18	73675	AAACAACAACAACA
19	72783	TTTTGTTGTTGTTT
20	72783	TGTTGTTGTTGTTT
21	72518	TGTGTGTATGTGTGT
22	70994	ACACACATACACACA
23	68584	TGTGTATGTGTGTGT
24	68584	TGTGTGTGTATGTGT
25	68353	GTGTGTGTGTGTGTT
26	68353	TTGTGTGTGTGTGTG
27	68235	CACACACACACACAA
28	68235	AACACACACACACAC
29	67362	ACACATACACACACA
30	67362	ACACACATACACACA

マイクロサテライト以外で、非常に多いパターンがどのような配列が確認した。

マイクロサテライト以外で多い配列は、マイクロサテライトから繰り返しパターンがすこし変形したものであった。

次のような可能性が考えられる。

マイクロサテライトの読み間違いの可能性。

マイクロサテライトのつなぎ目という可能性。

このようなパターンとマイクロサテライトだけから構成される測定断片は、対象外とする。

4. 5. マッピングアルゴリズムのトライアル②

マッピングアルゴリズムの基本部分のみを作成し、全測定断片を用いて、課題の洗い出しを行った。

対象

ペアエンドの一方だけを使用

Liver100%サンプル

対象配列数

11,795,262

アダプター配列が見つかった配列&15塩基以上

対象配列数

8,764,801

マイクロサテライト配列除外数

2,461

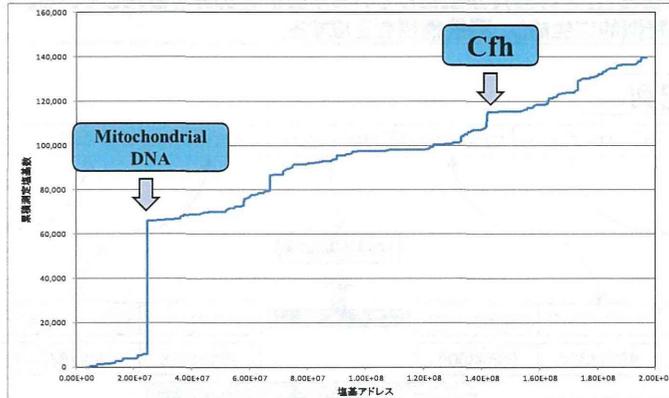
マッピングできた配列

4,591,681

マッピングできたペア(累計)

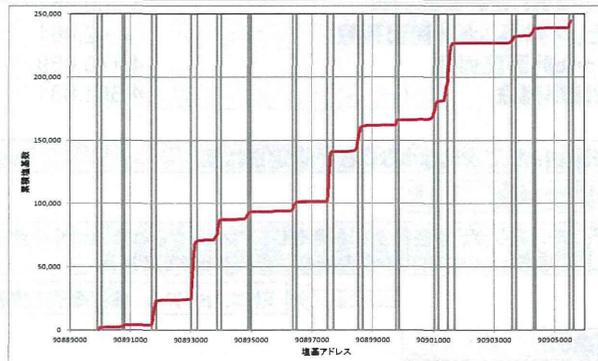
16,750,719

マッピングした配列を、塩基アドレスと、マッピング累積塩基数の積み上げ数でグラフを作成した。



多く合致している箇所の配列を調べたところ、染色体上の遺伝子ではなく、ミトコンドリアのDNAと一致した。

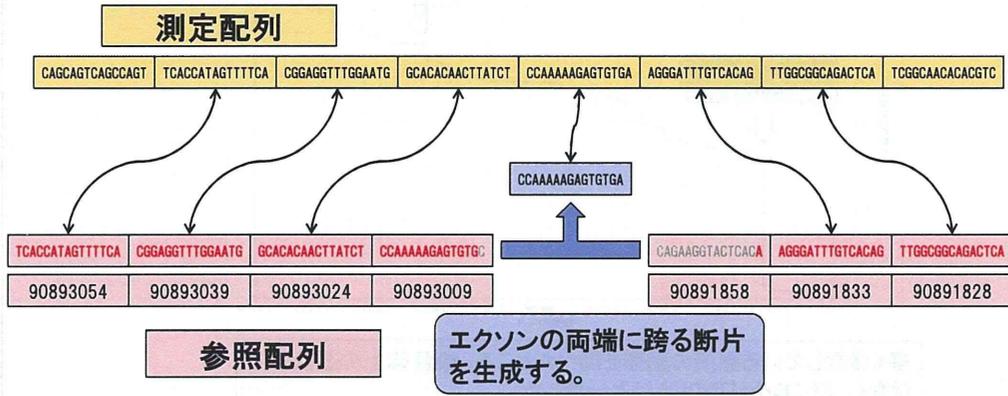
Alb近傍のマッピング断片約10万中、2千で累積分布を作成した。
簡易的に、断片中最も稀な配列(15塩基)位置で積み上げた。



エクソンで、積み上がっていく様子がうかがえる。
一部、イントロン領域にマッピングされている断片も存在した(スプライシング前のRNAと思われる)。

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。



対象計測配列数	11,795,262	
アダプターミスとして除外した計測配列数	3,030,461	25.69%
サイクロサテライトとして除外した計測配列数	2,461	0.02%
マッピングできなかった計測配列数	4,170,659	35.36%
マッピングできた計測配列数	4,591,681	38.93%

マッピングできなかった理由として次のようなことが想定される。

理由① ミトコンドリア由来

ミトコンドリアの配列が、染色体上にも存在し、マッピングされた。しかし、全ての配列が染色体上に存在しない。マッピングの漏れが存在するはずである。

ミトコンドリアも、参照配列とする。

理由② 読み間違い

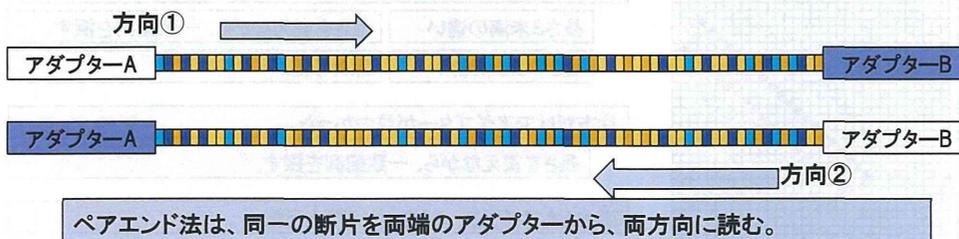
次世代シーケンサは読み間違いを起こす可能性がある。ペアエンド法は、同じ断片を別方向から読んでいて、精度を向上させる可能性がある。

ペアエンド法を活用して品質を向上させる。

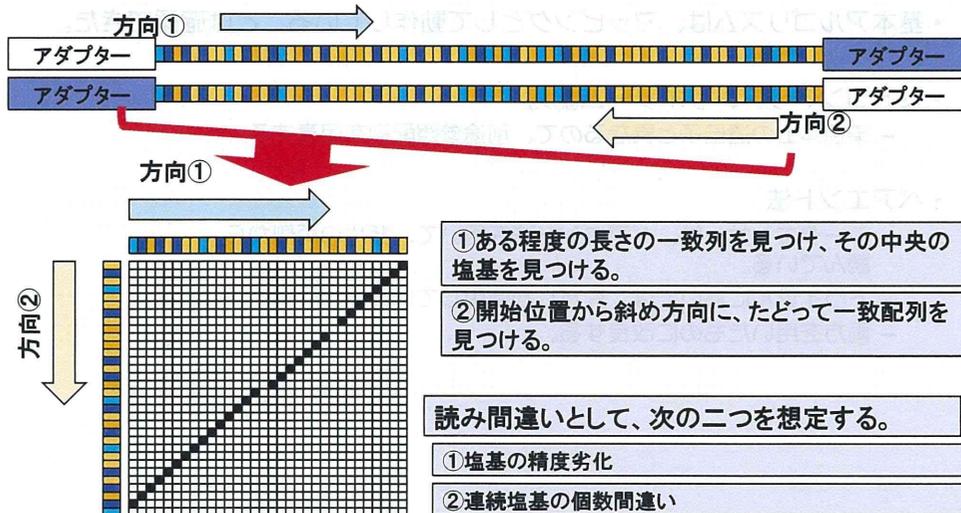
- 基本アルゴリズムは、マッピングとして動作していることは確認できた。
- ミトコンドリア、リボソーム配列
 - 染色体上の遺伝子と異なるので、別途参照配列を用意する。
- ペアエンド法
 - シークエンサーは、ペアエンド法を用いて、断片の両側から読んでいる。
 - トライアルにおいては、片方しか使用していない。
 - 両方を用いたものに改良する。

目的 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。

実装方法 ペア両方のファイルを読み込んで、処理しながらTeradataにロードする。

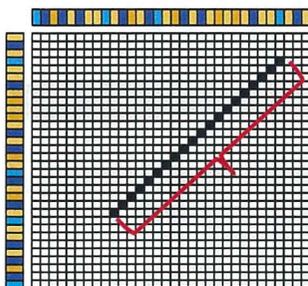


二方向からの配列から一致する箇所を見つけ出し、その状況により、適切な断片を求める。



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。
 一定のレベル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

除外する。

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。