

NCBI塩基配列データベースに蓄積するエラー解析

理化学研究所 バイオリソース研究センター 遺伝子材料開発室

三輪 佳宏、飯田 哲史、野崎 晋五、木嶋 順子、岸川 昭太郎、中島 謙一、
笹沼 俊一、大波 純一、中村 宣篤、村田 武英



問題提起: データベース精度を向上する方法

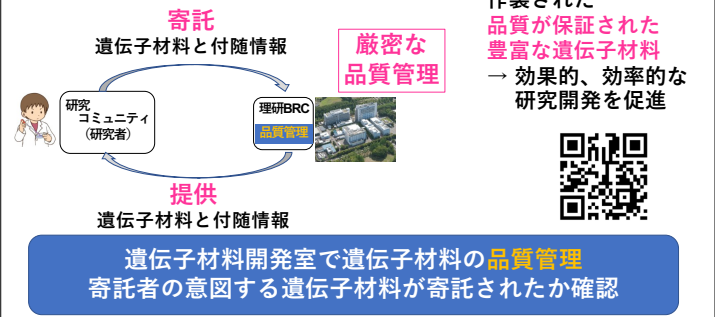
データベースにエラーが蓄積するには複数の原因がある

- 1) 実験そのものの精度の低さ
- 2) データのコピー&ペースト
- 3) 人為的ミス
- 4) バリエーション (エラーではない可能性)

原因に応じて、エラー防止のために求められる対策は全く異なる。
→実験研究者から情報研究者まで連携した対策づくりが必要

理化学研究所バイオリソース研究センター 遺伝子材料開発室

日本の遺伝子バンク



国際的塩基配列データベースのデータ検証

高度化した配列解析力+ 遺伝子バンクとしてのDNA資料

正しい塩基配列を確定させて、データベース上の配列エラーを探索・修正することが可能



検証結果の広報 「配列探偵」

1) 生物系月刊誌への連載

配列探偵 志久延子と弟子入り修行中の江良正、そして、配列探偵の後輩で理研BRCで研究員として働く穴理志津、の3人が、配列エラーに苦しむ研究者を救うためにデータベースの闇と戦う物語 紹介データは全てリアル!



2) RIKEN NEWS

NCBIの塩基配列データベースを検証し、蓄積したエラーを探し出すリアル配列探偵達の紹介

RIKEN NEWS 「データベースに潜む塩基配列エラーに警戒せよ」



塩基配列エラー実例1 puromycin耐性遺伝子

Puromycin-resistant gene(s?)

```

1 V E C P K D R A T W C M T R K P G A *
  .TGCGAG gTcCcgaaGgaccGcg ACCTGGTGCATGACCCGCAAGCCCGTGCTGA
2 .TGCGAG gTcCcgaaGgaccGcg ACCTGGTGCATGACCCGCAAGCCCGTGCTGA
  V E V P E G P R T W C M T R K P G A *
  
```

品質管理検査
2の配列しか出てこない

REFERENCE 1 (bases 1 to 806)
AUTHORS Lacaille,R.A., Pulido,D., Vara,J., Zalacain,M. and Jimenez,A.
TITLE Molecular analysis of the pac gene encoding a puromycin N-acetyl transferase from *Stenotrophomonas maltophilia* strain alboniger
JOURNAL Gene 79 (2), 375-380 (1989) **登録**

REFERENCE 2 (bases 1 to 806)
AUTHORS Jimenez,A.
TITLE Direct Submission
JOURNAL Submitted (17-MAR-1989) Antonio Jimenez, Centro de Biologia Molecular Severo Ochoa (CSIC/UAM), Universidad Autonoma de Madrid, Madrid, 28049, Spain **コピー**

COMMENT GenBank staff is unable to verify source organism and sequence and/or annotation provided by the submitter.
On Apr 11, 1995 this sequence version replaced gi: **修正**

2022年に新規登録

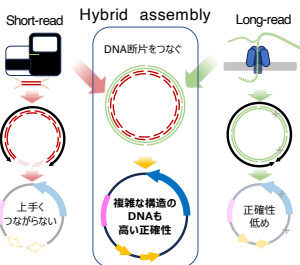
Puromycin-resistant gene in Addgene

Short readのみ → Assembleミス

← 実在しない幻のタンDEMリピーT

塩基配列の確定法

HTSを用いたDNA配列検査

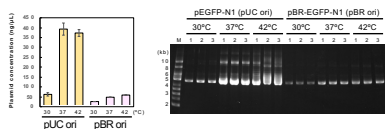


Short readシーケンサーだけでは、タンDEMリピーTなどの繰り返し配列を正しく解析できない。(左)

Long readシーケンサーだけでは、細かな塩基配列の正確性や連続塩基などが正しく解析できない。(右)

塩基配列エラー実例2 pUC-ori

..CGGTACTACTAGAAGAACAGTATTTGGTATCTGCGC.. 正しいpUC oriの配列
|||||
..CGGTACTACTAGAAGGACAGTATTTGGTATCTGCGC.. データベース上の配列



データベースと同じ配列の変異を導入 (pBR) して大腸菌でのコピー数を解析した結果 1塩基の違いでコピー数は激減する

MCS以外は同じ塩基配列の“兄弟ベクター”である pUC18とpUC19のうち pUC19のみが2000年に修正されたがpUC18はエラーのまま未修正

エラーのまま未修正

塩基配列エラー実例3 EF-1α プロモーター

「EF-1α promoter」という1つの名前でも異なる塩基配列が流通してしまっているが、そのことを認識している研究者は少ない

4 types of EF-1 alpha promoter

Type	genomic sequence	clones
Type 1	RDB18336 : pAY5 (mAID-EGFP-NLS piggyBac)	1
Type 2	Dr. Nagata JBC paper	11
Type 3	Dr. Nagata pEF-BOS plasmid	158
Type 4	lentiv. modified?	34