

特徴的塩基配列とステム構造に基づいた

DNA からの snoRNA 遺伝子検出法

薛晨¹ 山森一人² 剣持直哉³ 吉原郁夫²

宮崎大学工学研究科¹ 宮崎大学工学部² 宮崎大学フロンティア科学実験総合センター³

1. はじめに

タンパク質への翻訳情報を持たない非コード RNA と呼ばれる機能性 RNA の発見と役割の解析は、分子細胞生物学とバイオインフォマティクス双方において、現在最も重要な研究課題の1つである。基本的な代謝から個体発生、細胞分化までの様々な生命現象に関与する機能性 RNA が現在までに数多く見出されている[1]。機能性 RNA と疾患との関わりに関する研究成果も次々に報告されており[2]、創薬や再生医療分野などで大きな進展をもたらすことが期待されている。本稿では、機能性 RNA の1つである核小体低分子 RNA (snoRNA) をコンピュータによって自動検出する方法を提案する。

snoRNA は核小体内に存在し、リボソーム RNA 前駆体のスプライシングに関与する。snoRNA は図1に示したその構造から、Box C/D 型と Box H/ACA 型の2つに分類される[3]。2種類の snoRNA はそれぞれ独自の特徴的塩基配列を持つ。1つはボックスと呼ばれ、Box C/D 型 snoRNA は Box C (AUGAUGA)、Box D (CUGA)、Box C' (UGAUGA) と Box D' (CUGA) を持つ。Box H/ACA 型 snoRNA は Box H (ANANNA、N は任意塩基) と Box ACA (ACA) を持つ。もう一つは相補対(A-U と C-G)により形成されるステムと呼ばれる構造である。

提案手法では、特徴的塩基配列とステムの存在確率を評価することで snoRNA 遺伝子の自動検出を行う。

2. 提案手法

2.1. 提案手法の概要

snoRNA 遺伝子検出にあたり、まず対象とする塩基配列にボックスが含まれる確率を調べる。次に、ボックスが一定以上の確率で含まれる塩基配列のしかるべき位置に、相補対により形成されるステム構造が存在する確率を計算する。2つの確率の平均値がしきい値を超える場合、対象とする塩基配列に snoRNA 遺伝子が含まれる

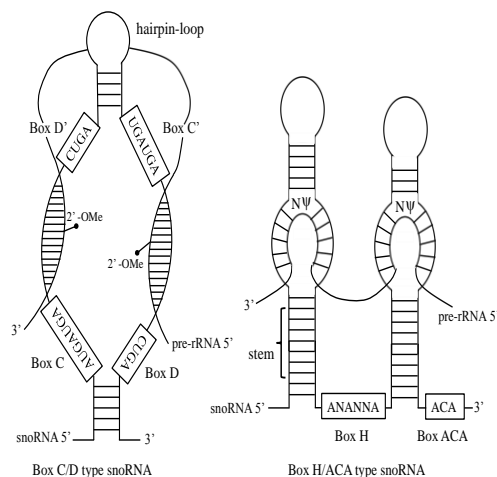


図1. Box C/D 型 snoRNA と Box H/ACA 型 snoRNA の構造

と判断する。

2.2. ボックス存在確率の導出

ある塩基配列 X 「ACCUGAU」の中に Box D が含まれるか否かを判定する場合について説明する。Box D は4塩基からなるので、X の先頭4塩基と Box D 配列をまず比較する。この場合、すべての塩基座において塩基が一致しないので、ボックスを構成する塩基配列と一致する割合を表す類似度は 0%となる。次に、切り出す領域を一塩基右にずらして「CCUG」と Box D を比較し、類似度 25%を得る。以下、X の末端となるまで同様の処理を行い、類似度の最大値を X における Box D の存在確率とする。この時、存在確率が 30%以下のボックスは存在しないものとして扱う。他のボックスについても同様の処理を行い、それぞれの存在確率を計算する。

2.3. ステム存在確率の導出

ある塩基配列 Y [...(Box D')GUACGUAA(Box C')...]におけるステム存在確率を求める場合について考える。ボックス間の配列

「GUACGUAA」の両端から内側に向かって、ステムを構成しうる相補対を探索する。8塩基からなるボックス間塩基配列がステムを構成すると仮定した場合、相補対数は4となる。実際には、(U-A)が1つ、(A-U)が1つ、(C-G)が1つの3つが含まれる。ステム存在確率は、構成しうる相補対数に対する、実際の相補対数として定義する。つまり、この例でのステム存在確率は75%となる。

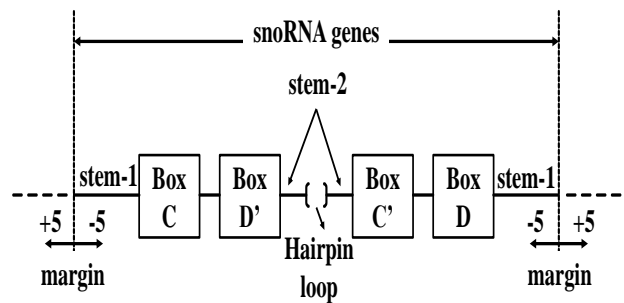


図 2. Box C/D 型 snoRNA 二次構造の展開図

3. 実験と考察

3.1. 実験データ

データベース「snOPY」[4]から取得した、長さ120のsnoRNA遺伝子全体を含む塩基配列1000個を正例として準備する。負例は正例と同じ生物種、同じ長さの、snoRNA遺伝子から離れた位置にある塩基配列1000個とした。評価スコアとして、正例と負例それぞれを正しく判定できる確率である識別精度、正例を正しく判定できる確率である感度、負例を正しく判定できる確率である特異度を用いる。

3.2. 実験結果と考察

予備実験により、ボックス存在確率とステム存在確率の平均が73%以上であればsnoRNA遺伝子ありとして判定を行うこととした。正例1000個と負例1000個で実験を行ったときの識別精度、感度、特異度をそれぞれ表1に示す。

表 1. snoRNA 遺伝子の検出結果

型	識別精度	感度	特異度
Box C/D	0.856	0.877	0.835
Box H/ACA	0.547	0.407	0.687

表1の通り、Box C/D型snoRNA遺伝子では約86%という識別精度を得ることができた。しかし、Box H/ACA型snoRNA遺伝子では約55%という識別精度しか得られなかった。これは、Box H/ACA型snoRNA遺伝子が持つBox H内に任意塩基を表すNが多く含まれており、その特徴を捉えにくいと考えられる。

3.3. Box C/D型snoRNA遺伝子検出実験

Box C/D型snoRNAの検出精度を評価するため、長さ5000の塩基データから正しくsnoRNA遺伝子の位置を特定できるかどうかの実験を行った。

データベース「snOPY」からsnoRNA遺伝子を含む5000塩基長のデータ100個を切り出し、あらかじめsnoRNA遺伝子が何塩基目から何塩基目まで占めるか記録しておく。次に、提案手法で

snoRNA遺伝子の有無を調べ、snoRNA遺伝子が含まれていた場合はその位置を求める。これらの位置を比較し、検出されたsnoRNA遺伝子の位置が正しいかどうか評価する。

図2に示すように、両端のstem-1から±5塩基の検出誤差を許して実験を行った。その結果、87個のデータでsnoRNA遺伝子を正しく検出することができた。このことより、提案手法の有効性が示された。

4. おわりに

本研究では、特徴的塩基配列とステム構造に着目した、DNAからのsnoRNA遺伝子検出法を提案した。実際の塩基配列を用いて識別精度を検証した結果、Box C/D型snoRNAについて約85.6%の識別精度を得ることができた。

今後の課題としては、Box H/ACA型snoRNAをより精度良く検出するため、snoRNA遺伝子の修飾情報を加味した手法の開発などが挙げられる。

参考文献

- [1] Eddy S (2001), Noncoding RNA genes and the modern RNA world. *Nature Reviews Genetics*. 2, pp. 919-929
- [2] 剣持直哉 (2004), リボソームと疾患. *実験医学 増刊*, 22, pp. 200-205
- [3] 剣持直哉 (2006), ゲノム情報に基づくsnoRNAの解析: ゼブラフィッシュを用いたアプローチ. *機能性 Non-coding RNA*, pp. 83-98
- [4] snoRNA Orthological Gene Database: http://snopy.med.miyazaki-u.ac.jp/snorna_db.cgi

A method to detect intronic snoRNA genes using characteristic nucleotide sequences and stems structures

¹Graduate School of Eng., Univ. of Miyazaki

²Faculty of Eng., Univ. of Miyazaki

³Frontier Science Research Center, Univ. of Miyazaki