複数細胞種における核酸結合タンパク質の 発現制御遺伝子と機能を予測する手法の開発

大里直樹 佐藤健吾

東京科学大学



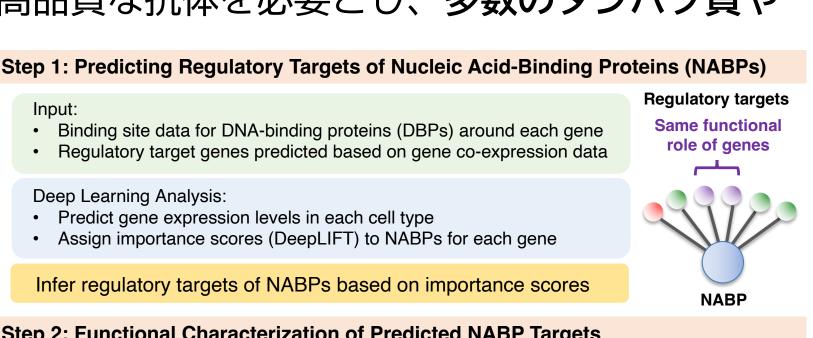


1. 要旨

- 核酸(DNA/RNA)結合タンパク質(Nucleic acid-binding proteins; NABPs)は、細 胞種特異的な制御や機能に関わるが、その**制御標的遺伝子や生物学的役割**は十分に解明 されていない。
- 遺伝子の共発現相関を利用した深層学習の手法を開発し、NABPの結合位置やモチーフ の情報を用いずに、NABPの標的遺伝子を予測し、その機能を推定した。
- 共発現に基づく標的遺伝子を用いることにより、遺伝子発現予測の精度が向上した。
- 予測されたNABPの標的遺伝子は、ChIP-seqとeCLIPの実験による標的遺伝子と高い整合 **性を示**し、ランダムに選んだ遺伝子よりも優れた性能を示した。また標的遺伝子を含む 共発現遺伝子は、DNA結合タンパク質のノックアウト実験結果と高い整合性を示した。
- ・ 機能アノテーション解析およびChatGPTを用いた機能推定により、AKAP8によるがん細 胞での概日リズム制御や、PKMによる解糖系制御など、**生物学的な意味のある機能**が同 定された。
- 深層学習・共発現ネットワーク・大規模言語モデルを組み合わせることで、**既知および** 新規のNABP機能を細胞種特異的に体系的に解析できる。

2. 背景

- 核酸結合タンパク質は、DNA結合タンパク質(DBPs)やRNA結合タンパク質(RBPs を含み、**細胞種特異的な制御や機能**があるが、網羅的な解析は困難である。
- ChIP-seq や eCLIP などの実験解析は高品質な抗体を必要とし、多数のタンパク質や 細胞種に適用することは難しい。
- RBP においては、明確なRNA結合モ チーフ配列を示さない場合があり、 またディスオーダー領域(IDR: Intrinsically Disordered Regions) を介した相互作用が知られており、 標的遺伝子の同定が困難である。
- 本研究では、遺伝子共発現データを 用いて、細胞種ごとのNABPの標的遺 伝子と機能を予測する。



Step 2: Functional Characterization of Predicted NABP Targets

- Perform enrichment analysis (e.g., GO, PANTHER Protein Class, Pathways, Reactome)
- Apply a ChatGPT-based method to infer functions and retrieve supporting literature Integrate ontology-based and semantic approaches
- Identify functional roles of NABPs in a cell type–specific context

3. 方法

- 深層学習を用いて、複数のヒト細胞種(例:HFF, HMEC, NPC, HepG2, K562)におけ るDNA結合タンパク質のDNA結合位置のデータから遺伝子発現量を予測した。
- DNA結合位置は、転写開始点(TSS)から±1 Mb以内の遠位エンハンサーを含み、 ChIP-seqデータ(GTRD)により同定され、eQTLデータにより遺伝子と対応づけされた。
- プロモーターおよびエンハンサーを統合し、1,310種類のDBPsに対応する230ビンの 入力マトリクスを構築し、DeepLIFTを用いたContribution scoreの解析を行った。
- Contribution scoreが低いDBPsを、共発現データベース(COXPRESdb)から取得した 核酸(DNA/RNA)結合タンパク質(NABPs)に置き換えた。
- 共発現は Pearson相関係数の zスコア > 2 を選び、共発現遺伝子のTSSから -2 kbまた は 0 bpの位置にNABPsの結合を仮定した。
- 上位あるいは下位にランクされたNABP-遺伝子ペアを選び、その標的遺伝子群について PANTHERにより機能エンリッチメント解析を行なった。
- さらに、ChatGPTを活用した手法により、遺伝子セットの機能推定および関連文献の抽 出を行なった。手動による文献検索も行った。
- 実験データによる検証として、DBPの場合はChIP-seq、RBPの場合はeCLIPデータ(ENCODEおよびGEO)を用い、高スコア遺伝子との一致を評価した。DBPのノックアウ ト実験により発現量が変動した遺伝子と共発現遺伝子を比較した。

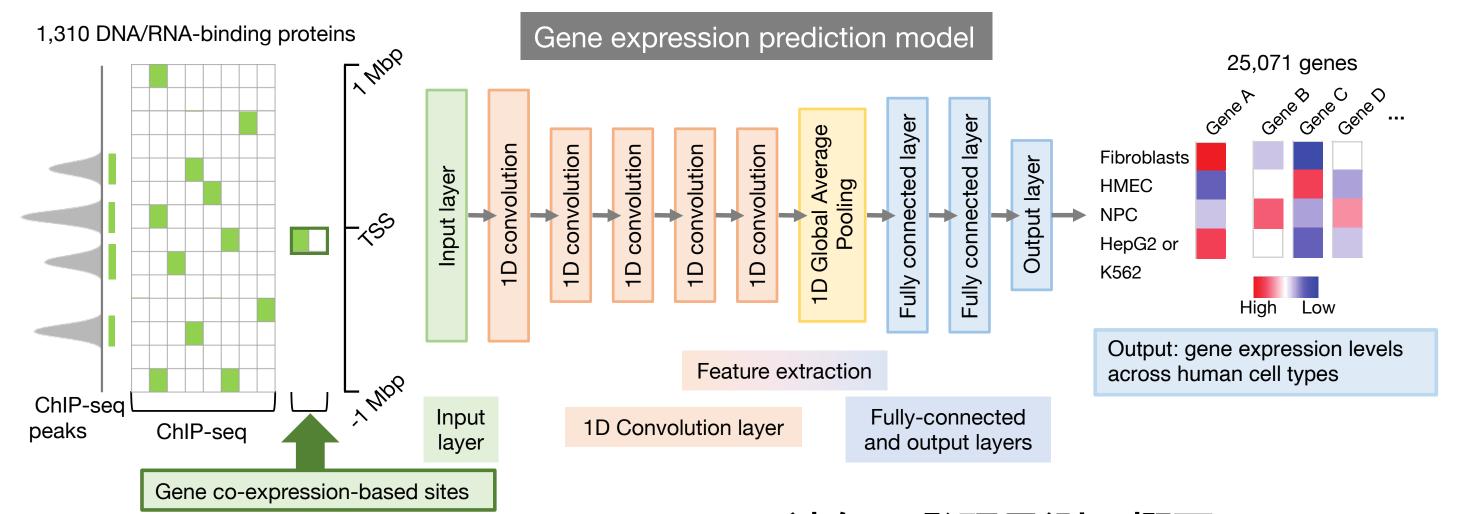


Figure 1. 遺伝子発現予測の概要

DNA結合タンパク質のゲノム上の結合位置(プロモ ーターと遺伝子遠位)の情報を用いて、複数のヒト 細胞種の遺伝子発現量を予測した。

Figure 2. 遺伝子の共発現(Gene Coexpression)

遺伝子の共発現は、2つの遺伝子間の発現プロ ファイルの類似性を示す。

機能的に関連する遺伝子や同じ調節機構により 制御されている遺伝子は、一般に細胞や組織間 で同様に発現変動する傾向がある。

ある核酸結合タンパク質(NABP)をコードす る遺伝子が他の遺伝子と類似した発現パターン を示す場合、遺伝子発現が核酸結合タンパク質 により制御され、機能的に関連する可能性があ る。

4. 結果

- Contribution scoreの低いDNA結合タンパク質(DBPs)を遺伝子共発現データに置き 換えると、遺伝子発現予測の精度が向上した(DBPの場合:R = 0.70 → 0.80、RBPの 場合:R = $0.70 \rightarrow 0.81$) (Figures 1 and 2)。
- NABPに対してChIP-seqまたはeCLIPで実験的に同定された結合位置をもつ遺伝子は、 Contribution scoreの上位または下位にランクされる傾向を示し、遺伝子発現制御にお いてより強い活性または抑制の役割が示唆された (Figure 3a)。
- Contribution scoreの高い遺伝子は、ランダムな遺伝子よりも**高い割合でChIP-seqおよ** びeCLIPピークと共局在することが確認された(Figure 3b)。
- 遺伝子共発現データに基づくDBPの標的遺伝子のほぼすべてが、DBPのノックアウト実 験により発現量が有意に変動した遺伝子に含まれた。
- 10種類のタンパク質について予測されたNABPの制御標的遺伝子の機能エンリッチメント解析 では、UniProtのアノテーションと一致するGene Ontologyとの関連が示された (**Table 1**)。
- ChatGPTを用いた機能予測もこれらの役割と同様の結果を示し、さらに新規の関連性を 同定した(例:AKAP8による概日リズム制御、PKMによる解糖系制御)。これらの機 能は文献により確認できた(Table 1)。
- 深層学習の学習後には、概日リズム制御やキナーゼ活性に関連する機能アノテーション が特異的に得られ、学習前の遺伝子では見られなかった特徴が観察された(Table 1)。
- 共発現データに基づく深層学習とAIによる文献解析を組み合わせた本手法は、細胞種特 異的なNABPの制御的役割や機能を同定できる(Figures 4 and 5)。

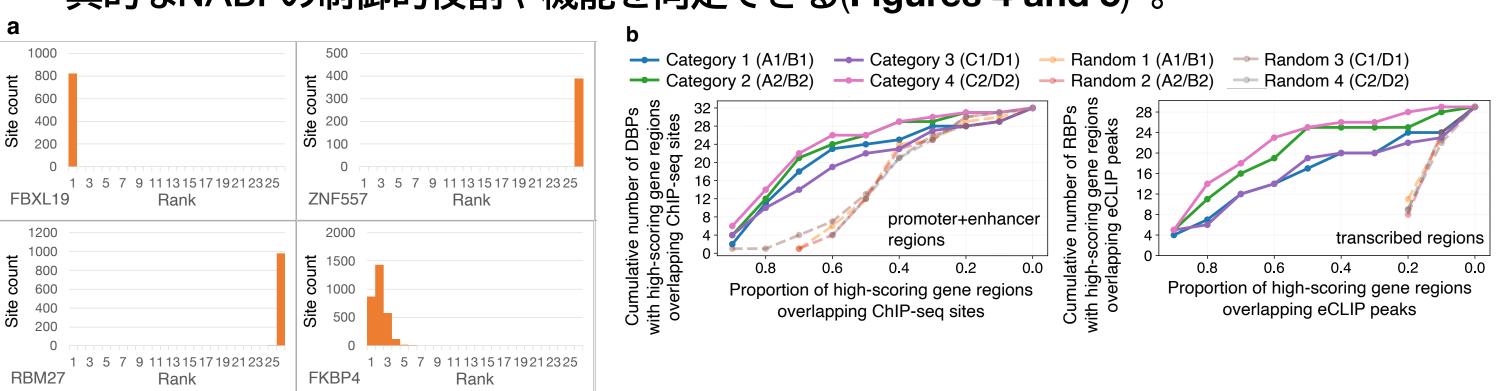


Figure 3. 予測された制御標的遺伝子とChIP-seqおよびeCLIP結合データの比較

C2H2 zinc finger transcription factor negative regulation of RNA biosynthetic DNA-binding transcription activator activity Circadian clock system RNA polymerase II-specific RNA polymerase II transcription regulatory region sequence-specific DNA binding Circadian-regulated energy metabolism and insulin signaling (0.74) Transcriptional regulation by C2H2 zinc finger proteins and signaling modulation (0.72) glycolytic process ATP metabolic process chromatin/chromatin-binding, or -regulatory Hypoxia-adaptive metabolic remodeling transcription regulator activity and cell migration (0.79) metal ion binding RNA-binding protein-mediated transcriptional integration and pre-RNA polymerase II transcription regulatory mRNA processing (0.89) region sequence-specific DNA binding **ZKSCAN5** DNA-binding transcription factor activity Ribosome biogenesis and RNA C2H2 zinc finger transcription factor and nuclear export (0.91) RNA polymerase II cis-regulatory region sequence-specific DNA binding metal ion binding large ribosomal subunit RNA processing, ribonucleoprotein assembly, and transcription-coupled ribonucleoprotein complex quality control (0.83) cytosolic large ribosomal subunit Mitochondrial-linked ribosome function and RNA-guided surveillance in **ZNF446** translational homeostasis (0.84) C2H2 zinc finger transcription factor Transcription-coupled RNA processing and mitochondrial RNA metabolism chromatin binding transcription coregulator activity **ZNF557** Chromatin-associated transcriptiona DNA-binding transcription activator activity, RNA polymerase II-specific RNA polymerase II transcription regulatory region sequence-specific DNA binding K562 cells ion binding, metal ion binding (cancer cells) neuronal cells Transcriptional regulation and RNA metabolic coordination via zinc finger proteins and RNA-associated complexes

Table 1. 核酸結合タンパク質の制御標的 遺伝子の機能注釈と機能推定

核酸結合タンパク質(NABP)の機能アノテ ーションは、UniProtデータベースから取得し た。さらに、予測された制御標的遺伝子にお いて有意にエンリッチされたアノテーション について、PANTHERデータベースを用いて同 定した。各NABPとその予測標的遺伝子で共 通または類似の機能アノテーションを表に示

ChatGPTベースの手法により推定された機能 的役割は黄色でハイライトされている。括弧 内の数値は信頼スコア(confidence score) を示し、対応する生物学的プロセスに関与す る標的遺伝子の割合を反映している。信頼ス コアの範囲は0.00~1.00であり、高(High : 0.87-1.00、中(Medium): 0.82-0.86、低(Low):0.01-0.81、なし(None):0 として分類される。).

Figure 4. 核酸結合タンパク質の細胞種に特 異的な制御と機能

K562のがん細胞では、AKAP8がPER2やCRY2な ど複数の概日時計遺伝子の発現制御に関わり、代 謝リズムの制御への関与が示唆された。一方で、 神経細胞ではこのような関連はなく、AKAP8の細 胞種特異的な機能的役割が示された。

Figure 5. がん細胞のAKAP8とPKMによる低酸素応答性の遺伝子発現制御

PPP / NADPH

5. 結論

Hypoxia Adaptation / Tumor Survival

Circadian genes

R-loop formation / resolution

- 本研究は、遺伝子の共発現データから、核酸結合タンパク質(NABP)の制御標的および 機能を予測したことを示した。ChIP-seqやeCLIPと比較して、RNA-seq**は**高速・低コス トで、様々な細胞種や条件で利用可能である。
- 既知の遺伝子機能と整合し、文献的にも裏付けられる制御標的遺伝子の機能を新たに同 定した。
- 機能エンリッチメント解析およびChatGPTを用いた解析により、細胞周期制御や概日リ ズム制御などの既知および新規のNABP機能が明らかとなった。
- Contribution scoreにより、直接的および間接的な相互作用を予測でき、NABPによる制御 を包括的に理解できる。
- 細胞種特異的な予測を支援し、異なる生物学的コンテキストや複合体内で機能するNABP の発見につながると期待される。
- 本手法により、大規模なNABPの制御標的遺伝子と機能の予測ができ、新規の遺伝子発現 制御メカニズムの発見につながる。



ChIP-seq peaks of

A>G

associated SNP

DNA-binding proteins Gene expression level

eSNP: expression- TSS: Transcription start sites

eGene: expression-associated gene; a gene whose

States (Spatial, Temporal, Environmental)

PCC = 0.1

Gene C

expression level is associated with an eSNP

PCC = 0.7