

AI を用いたタンパク質間複合体予測から、

機能未知遺伝子の機能を推察する

河口 真一

大阪大学 大学院生命機能研究科

1. はじめに

細胞内で生命現象の物理化学的反応を担っているのは、多くの場合、タンパク質である。タンパク質は、20種類のアミノ酸が連なったポリマーであり、アミノ酸配列に応じた立体構造に折り畳まれる。生物のゲノム DNA にコードされているタンパク質の種類は、ヒトで約 20,000 程度、モデル生物のショウジョウバエでも、約 12,000 種類も存在する。世界中の研究者によって、各遺伝子がコードするタンパク質の機能を調べる研究が続けられており、多くの遺伝子について分子機能が明らかにされてきた。しかし、未だに機能が解明されていない遺伝子が 20%以上も残されており、分子生物学上の大きな謎となっている。これら機能未知な遺伝子の分子機能に関する手がかりを得る上で、直接的に相互作用する分子を同定することは、極めて有用である。なぜなら、機能がわかっているタンパク質との相互作用が明らかになれば、機能未知遺伝子も、同じ経路、あるいは関連する経路で働いている可能性が高いからである。相互作用するタンパク質を実験的に同定する研究も広く行われているが^{1,2}、時間と労力が多くかかるため、迅速なスクリーニング手法が望まれている。タンパク質間の相互作用は、各タンパク質の立体構造表面の特性によって決まる。そのため、タンパク質の立体構造情報を用いて、ドッキング法により結合を予測する試みも行われてきたが、未だ信頼できる予測法は確立されていない。

最近、DeepMind 社によって開発された AlphaFold2 という AI プログラムは、アミノ酸配

列の共進化情報をを利用して、タンパク質の立体構造を高精度で予測できるとして注目されている³。さらに、タンパク質の複合体構造も予測できることが期待されている⁴。そこで本研究では、AlphaFold2 を利用して、迅速なタンパク質複合体のスクリーニングを行い、機能予測の一助とする試みを試みた。

2. AlphaFold2 プログラムを用いたタンパク質複合体構造予測

AlphaFold2 による立体構造予測は、主に 2 つのステップに分けることができる。

- 1) アミノ酸配列の多重アライメント作成
- 2) 立体構造、及び複合体構造の予測

予測したいタンパク質のアミノ酸配列と類似したアミノ酸配列を、数多くの生物種から収集し、アライメントすることによって、互いに依存した変化（共進化）が推測される部分を抽出する。共進化している部分は、立体構造上、距離的に近接すると仮定して、立体構造のモデルを構築する。本研究では、ショウジョウバエのほぼ全てのタンパク質についての多重アライメントを作成・保存し、それらを再利用することによって、大規模な複合体予測に必要な時間を短縮させている。

AlphaFold2 は、数多くの立体構造を学習し、予測に用いるモデルパラメータを 5 パターン備えている。オリジナルの計算フローでは、5 つのパターンに対応する 5 つの構造を、1 つの GPU で逐一計算する仕様であったが、大阪大学サイバーメディアセンターの SQUID では、GPU を 5 つ用いて並列に計算するように変更し、予測時間をさ

らに短縮している（NEC ソリューションイノベータ株式会社との共同研究）。これらの変更によって、1 日、1 アカウントあたり、約 150 ペアのタンパク質複合体予測ができるようになり、通常のコンピュータと比べて、約 100 倍速く計算することが可能になった。

3. 20 種類のタンパク質間の複合体構造予測

3.1 信頼性スコアの分布

本研究では、ショウジョウバエの生殖細胞内に存在するタンパク質集積体の 1 つであるヌアージュ構造に着目した⁵。ヌアージュでは、piRNA という小分子 RNA が產生されており、piRNA と結合した PIWI ファミリータンパク質が、ゲノム DNA の損傷を引き起こすトランスポゾンの発現を抑制している。ヌアージュに局在する 18 種類のタンパク質と、機能的に関連するミトコンドリア局在タンパク質、2 種類を合わせた合計 20 種類のタンパク質について、1 : 1 の総当たりで 400 ペアの結合予測を行った。

表 1：本研究で用いた 20 種類のタンパク質

タンパク質	残基数	局在	ドメイン
vas	661	Nuage	DEAD
spn-E	1434	Nuage	DEAD
tej	559	Nuage	Lotus, Tudor
tapas	1222	Nuage	Lotus, Tudor
qin	1857	Nuage	RING, Tudor
Kots	892	Nuage	Tudor
krimp	746	Nuage	Tudor
squ	241	Nuage	
mael	462	Nuage	Meal
aub	866	Nuage	PAZ, Piwi
AGO3	867	Nuage	PAZ, Piwi
papi	576	Mitochondria	Tudor, KH
vret	691	Nuage	Tudor
bel	801	Nuage	DEAD
zuc	253	Mitochondria	PLD-like
cup	1117	Nuage	
tral	657	Nuage	Lsm
me31B	459	Nuage	DEAD
shu	455	Nuage	PPIase
BoYb	1059	Nuage	DEAD, Tudor

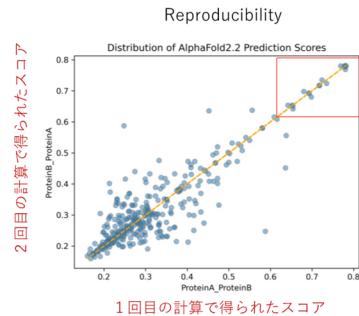


図 1：AlphaFold2 の予測スコア分布

AlphaFold2 は、各複合体の予測構造に対して、タンパク質単体での構造と、複合体界面の構造を評価し、信頼性スコアを出力する。スコアは、0 から 1 の範囲であり数値が大きいほど、信頼性が高いことを意味する。総当たりで計算すると、同じペアのタンパク質複合体を 2 回計算することになる。1 回目と 2 回目の信頼性スコアをプロットしたところ、スコアが低い（0.5 以下）場合には、スコアのばらつきが見られた。スコアが高い（0.6 以上）場合には、比較的再現性の良い結果が得られたことから、スコア 0.6 以上のペアを複合体形成の候補とした。

3.2 実験的検証

このスクリーニングにおいて、スコアが 0.6 以上のペアは、13 種類見られた。そのうち、7 つについては、既に結合が報告されているペアであった。このことは、AlphaFold2 が複合体ペアを予測できることを示唆している。残りの 6 ペアについては、これまでに報告がない新規な結合候補であった。その中で、Spn-E タンパク質と Squ タンパク質のペアに注目した。Spn-E は、RNA ヘリケースであり、ヌアージュにおける piRNA の产生に重要であることが知られている⁵。一方で、Squ タンパク質は、機能を推測できるドメインもなく、詳しい機能解析もされていない、機能未知タンパク質である。両者が複合体を形成するかどうか、実験的検証を行った（図 2）。

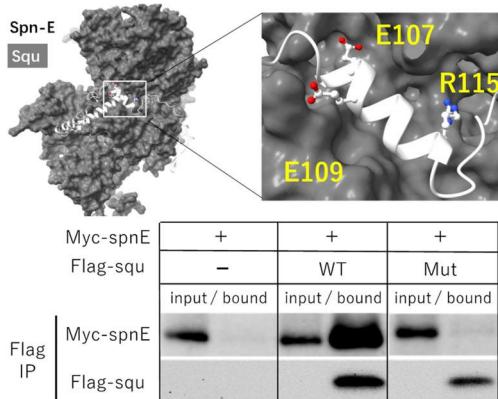


図2：AlphaFold2が予測した SpnE-Squ複合体と
共免疫沈降法による実験的検証

Mycタグ、Flagタグをそれぞれ付加したSpn-EとSquを昆虫培養細胞内で発現させた。Flagタグに特異的な抗体で免疫沈降を行い、Myc-Spn-Eの結合をMyc特異的抗体で検出した。その結果、Squタンパク質が共発現している場合にのみ、Spn-Eが結合画分に検出され、2つのタンパク質が結合し、複合体を形成することが明らかになった(図2下中央)。さらに、Squの α -ヘリックスの1つが、SpnE表面の溝に埋まっている構造が予測されたことから、この部分が結合に寄与していると考えられた。実際に、この α -ヘリックスとSpn-Eとの間で塩橋を形成すると予測された残基を変異させたところ、その変異Squは、Spn-Eとの結合が損なわれていた(図2下右)。したがって、AlphaFold2による予測構造は正しいことが示唆された。これらの結果からAlphaFold2を用いた複合体のスクリーニングが有効であると考えられる。機能未知であるSquの役割は、Spn-Eに結合することで、Spn-Eの安定性や活性を制御しているか、あるいは、Spn-Eが他の因子と相互作用することを制御していると推察される。

4. 細胞内の全タンパク質をスクリーニング

AlphaFold2による複合体スクリーニングを、細胞内の全タンパク質に対して適用することを試みた。ここでは、生殖細胞と体細胞の両方で、多様な機能が知られているPiwiタンパク質を用い

た。Piwiタンパク質は、piRNAと結合し、それと相補的な配列をもつトランスポゾンのmRNAを切断することによって、有害なトランスポゾンの発現を抑制している。さらに、核内では染色体のヘテロクロマチン化を誘導し、転写レベルで発現を抑制することも知られている⁶。

ショウジョウバエのゲノムには、約12,000種類のタンパク質がコードされている。Piwiタンパク質と全タンパク質との1:1結合スクリーニングを行った。大阪大学サイバーメディアセンターのSQUIDを用いて、1日に約100ペアの構造予測を行い、4ヶ月ほどで計算を終えることができた。その結果、大部分のペアについては、複合体に対する信頼性スコアが0.3未満と低く、結合していないと予測された。スコアが0.6以上のものは、全体の約1%に相当する177個見つかっている(図3)。

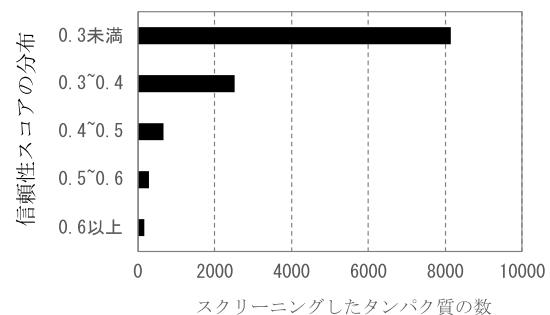


図3：Piwiとの複合体予測におけるスコア分布

この中には、これまでにPiwiと結合することが報告されているArxタンパク質に加えて、いくつかの機能未知タンパク質も含まれていた。特に、PiwiとCG33703という機能未知タンパク質の複合体については、高い信頼性スコア(0.82)をもつて予測された。CG33703遺伝子は、ゲノム上に多数のコピーがあることが知られているが、その機能については全く不明である。

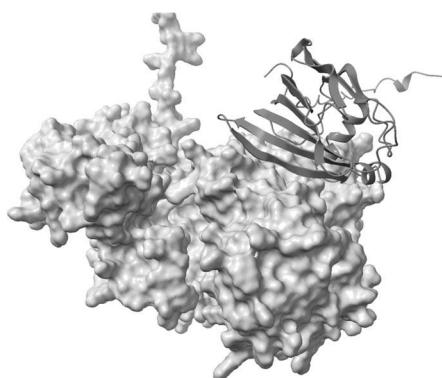


図4：予測された Piwi と CG33703 との複合体

Piwi: 空間充填モデル

CG33703: リボンモデル

今後、実験的検証を行い、新規な結合ペアを見出すことによって、機能解析を展開できると期待される。

5. おわりに

タンパク質のアミノ酸配列から立体構造を予測する試みは、古くから行われてきた。その過程で、2000年代に日本で行われたタンパク3000プロジェクトが、タンパク質の構造パターンの多くを明らかにしたことによって、ドメイン単位での立体構造予測が可能となった。現在のAlphaFold2では、ドメイン間の配置もかなりの精度で予測できるようになっている。

現在、AlphaFold2のようなAIプログラムの開発が盛んに行われており、生命科学の発展に大きく貢献することが期待されている。しかし、そのような計算機科学と、それを利用する実験科学の間には、まだ隔たりがあるのが実情である。両者の隔たりを埋めて、最新の計算機科学を、実験科学が迅速に取り込むことができる環境が望まれる。今後は、癌化した細胞などで過剰発現が見られるタンパク質や、実験的に扱うのが困難な毒性タンパク質の相互作用を *in silico* でスクリーニングし、新規な結合候補を示すことによって、タンパク質の機能解析を進展できると期待される。

このプロジェクトを進めるにあたり、サイバーメディアセンター応用情報システム研究部門の

伊達進教授、高性能計算・データ分析融合基盤協働研究所の曾我隆特任准教授（常勤）に大変お世話になりました。感謝の意を表します。

参考文献

1. Giot, L. *et al.* A protein interaction map of *Drosophila melanogaster*. *Science* **302**, 1727–1736 (2003).
2. Guruharsha, K. G. *et al.* A protein complex network of *Drosophila melanogaster*. *Cell* **147**, 690–703 (2011).
3. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
4. Evans, R. *et al.* Protein complex prediction with AlphaFold-Multimer. 2021.10.04.463034 Preprint at <https://doi.org/10.1101/2021.10.04.463034> (2021).
5. Lin, Y., Suyama, R., Kawaguchi, S., Iki, T. & Kai, T. Tejas functions as a core component in nuage assembly and precursor processing in *Drosophila* piRNA biogenesis. *J Cell Biol* **222**, e202303125 (2023).
6. Huang, X. & Wong, G. An old weapon with a new function: PIWI-interacting RNAs in neurodegenerative diseases. *Translational Neurodegeneration* **10**, 9 (2021).