

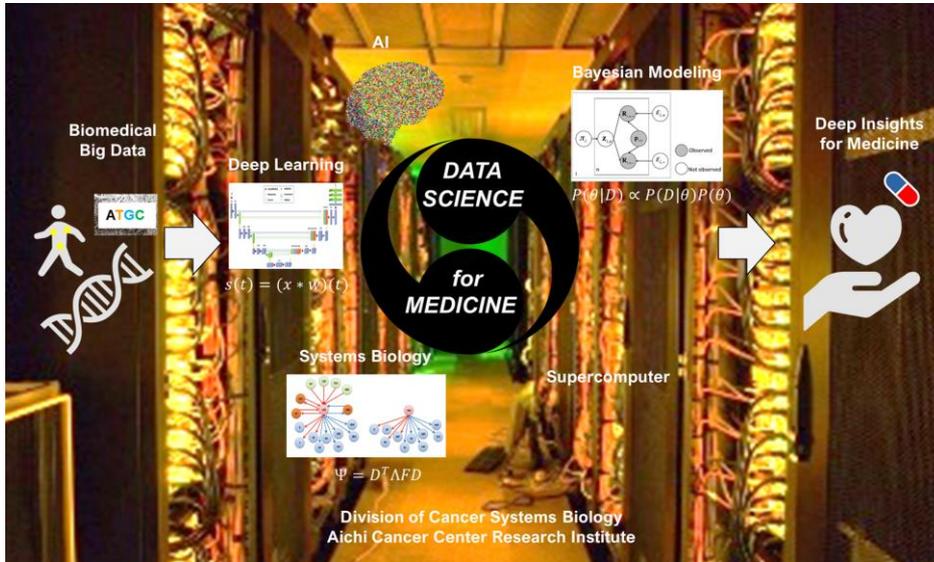
AIとスパコンを活用した がんビッグデータ 解析手法の開発

愛知県がんセンター
システム解析学分野



システム解析学分野 について

データ科学の力で医療へ貢献！



- がん細胞から得られたゲノムデータなどの生体ビッグデータをAIとスーパーコンピュータ（スパコン）を使って解析する方法の研究を行っています
- 患者さんのデータから、がん細胞の複雑なシステムに関わる情報を抽出し、一人ひとりに合わせた医療へつなげることを目指しています

研究室HP

<https://cancer-c.pref.aichi.jp/site/folder5/1154.html>



研究室メンバー

分野長： 山口 類

研究員： 郭 中梁

リサーチ

レジデント： 武藤 理

木曾田 暁

任意研修生： 末廣 智也

連携大学院生： 袁 一凡

研究補助： 鈴木 一基

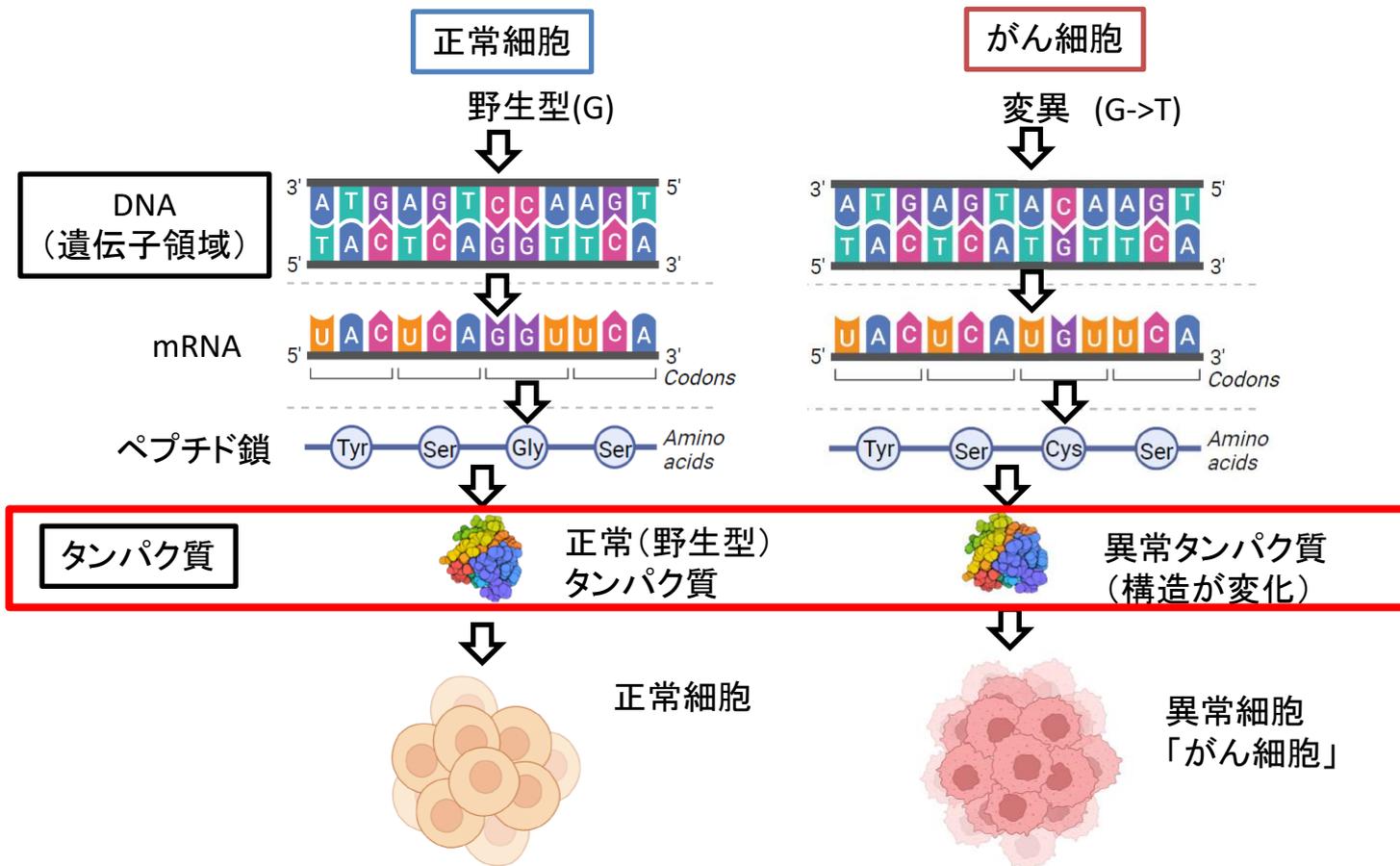
事務員： 竹中 亜由美



はじめに

- がん細胞は、正常な細胞中のDNA（ゲノム）配列に変異が起きることにより、タンパク質の構造が変化するなどして、細胞が異常に増殖する能力を獲得しています
- 計測技術の発展により、DNA等の生体分子を計測した大量のデータを得ることができるようになっています
- がん細胞中のDNAに変異がおきる場所やパターンは多様であり、がん細胞の性質を知ることや薬を選択するために、コンピュータを使って計測データからDNAの塩基配列を決定し、またどこにどのような変異が生じているかを正確に検出する必要があります（**研究1**）
- さらに変異の結果タンパク質の性質がどのように変化するかを予測することも重要です（**研究2**）
- 研究のサイクルを加速し成果を現場に還元するには、得られた大量のデータを迅速かつ簡単に解析できるシステムも必要です（**研究3**）
- ここでは、近年発展している、がんのビッグデータ解析の状況をシステム解析学分野での研究を交えて概説します

DNAと変異



細胞中の**DNA**には**遺伝子**と呼ばれる領域（約2万か所）が含まれています。**遺伝子領域**には**タンパク質の設計図**となる情報が含まれています。

多くの場合がん細胞では**遺伝子に変異**が入り、**構造が変化した異常なタンパク質**が作られることにより、がんの性質に影響を与えています。またその変異は多様です。そのため**遺伝子変異を正確に検出しタンパク質の性質を推測**することが重要です。

背景：DNA等の生体分子が大量かつ高速に計測可能に

次世代シーケンサー(NGS): DNA配列等を読む機械



2021年8月時点で約512ドル

近い将来100ドルを切る予想
(解析費用は含まない)

大規模ゲノム解析プロジェクト
Genomics England

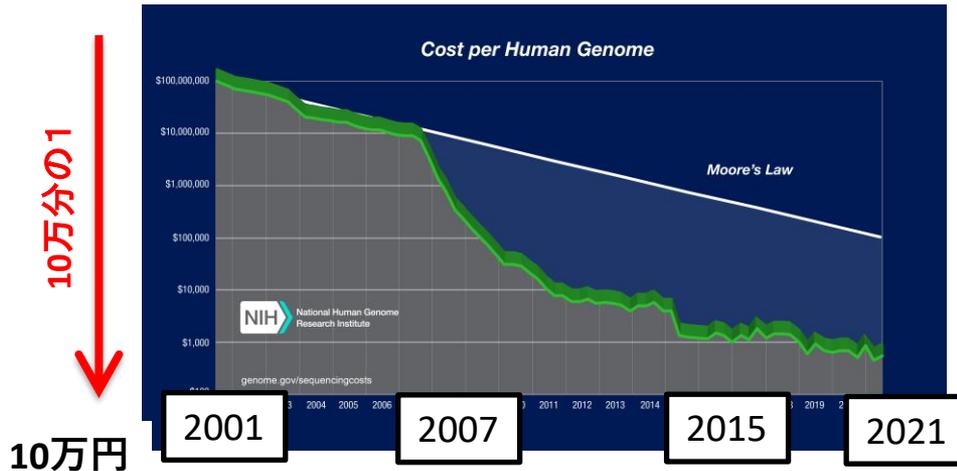
All of Us

ToMMo

厚労省・ゲノム解析実行計画
AMED 大規模臨床がん全ゲノム
データ解析プロジェクト

NHGRIの予測: 2025年までに
40ExaB (10⁹GB)のシーケンスデー
タが全世界で産生
全ゲノムに換算すると
約2億人分

100億円 一人分のゲノムを決定するためのシーケンスコスト



<https://www.genome.gov/about-genomics/fact-sheets>

次世代シーケンサー(NGS)と呼ばれる計測機械の性能向上により、一人分の全ゲノム情報を得るのにかかるコストは20年で10万分の1になりました。

これらのシーケンスデータからがん細胞特有の変異・変化を見つけることで、がんの原因や治療法に関する情報を得ることができます。

現在、保険診療で行われている**がんゲノム医療**でも役立てられています。

シークエンサーって、 どんなデータ？

生のサンプル

DNAとして抽出
ATCGの4種の文字



次世代シークエンサー



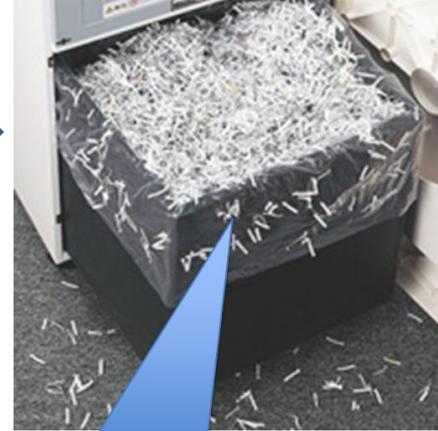
得られるデータ

ゲノムシュレッダー

100文字ぐらいの断片になった

21億ピース

の文字列断片がコンピュータに
吐き出される



ATCCGGTAAAT.....TTCA

100~150塩基

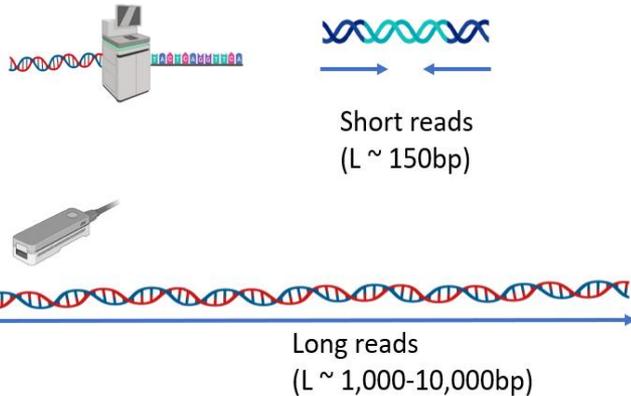
※全ゲノムシーケンスを想定

一人分のゲノムは、4種類の塩基(A、T、G、C)からなる塩基対、約30億個分からなります。現在主流のシークエンサーでは、一度にはDNA中の約100塩基対分しか読めません。これはATGCの4文字からなる30億文字分の文章が書かれた書類の束が、シュレッダーにかけられ100文字ずつの断片になって出てくるようなものです。通常、がん細胞特有の変異を見つける場合、がん細胞由来の書類を40コピー分、正常細胞由来の書類を30コピー分、シークエンサーで読み取ります。つまり21億ピースの文字列断片が得られます。ここから各ピースがゲノム上のどこにあったかを正確に決定し、がん細胞と正常細胞の違いを見つけ出す(変異を検出する)必要があります。

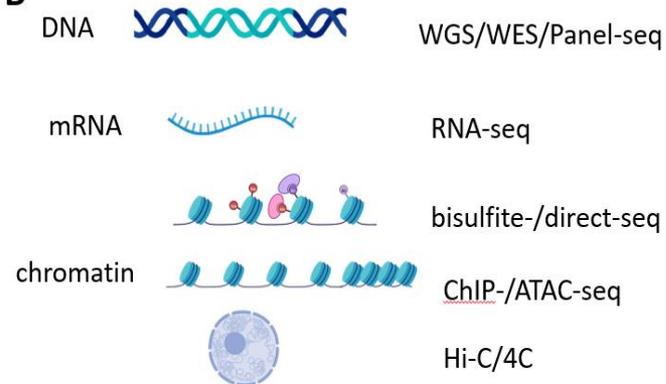
生体分子データの多様化

NGSでは、DNAに限らずRNA等他の分子の情報も得ることができます。またこれらの分子は、多数の細胞からなるがん組織だけでなく、単一の細胞ごとや腸内細菌叢から得ることもできるようになっています。このような複雑なデータの解析にはAIや機械学習（ML）の手法が必要です。

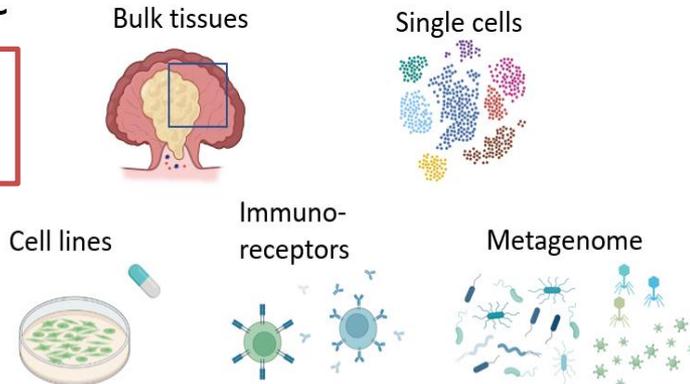
A



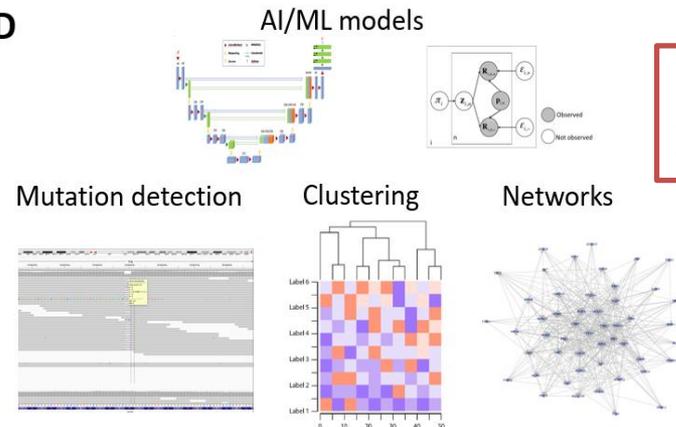
B



C



D

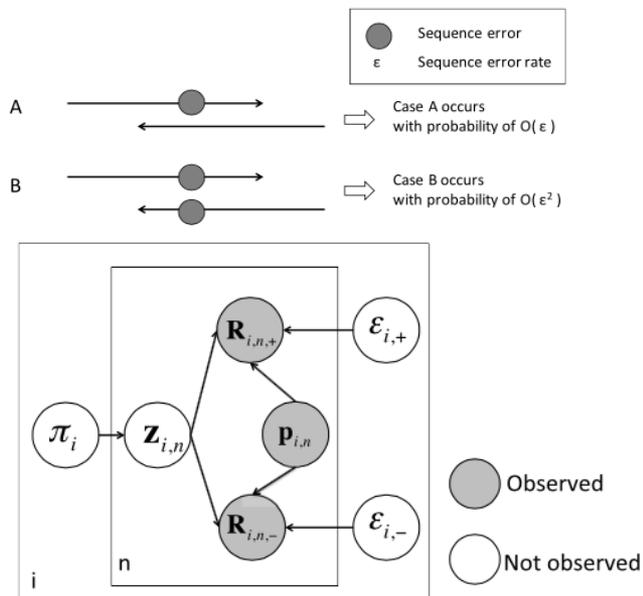


生体分子データから有用な情報を抽出する 手法の開発を行っています。

ベイズ統計モデル

×

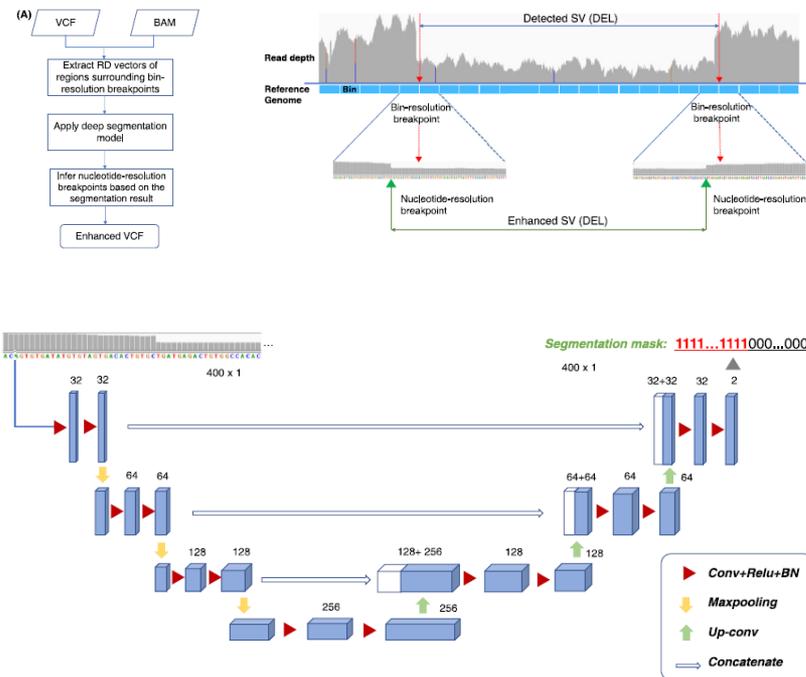
深層学習モデル (AI)



$$p(R_i, Z_i | \gamma_i, \alpha_{i,+}, \alpha_{i,-}) = p(\pi_i | \gamma_i) p(\epsilon_{i,+} | \alpha_{i,+}) p(\epsilon_{i,-} | \alpha_{i,-}) \cdot \prod_n p(R_{i,n,+}, R_{i,n,-} | z_{i,n}, \epsilon_{\pm,i}, \pi_{i,n}, p_{i,n}) p(z_{i,n} | \pi_i)$$

短い塩基変異の検出法

Moriyama et al., Bioinformatics, 2019

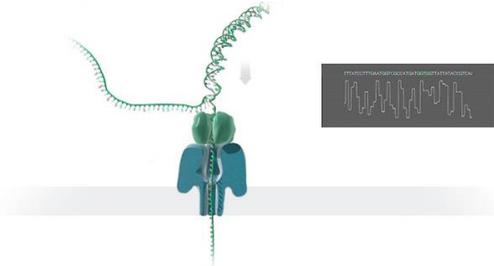


長い塩基変異の検出法

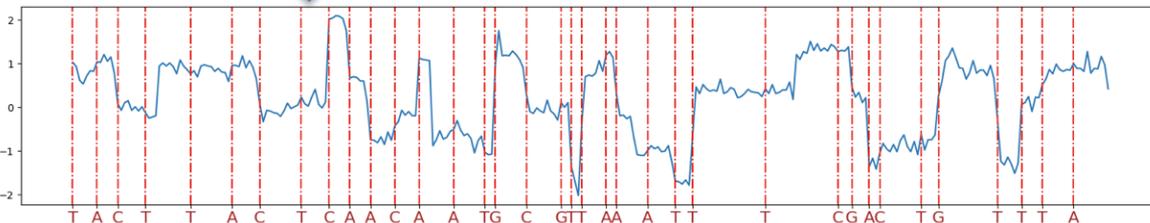
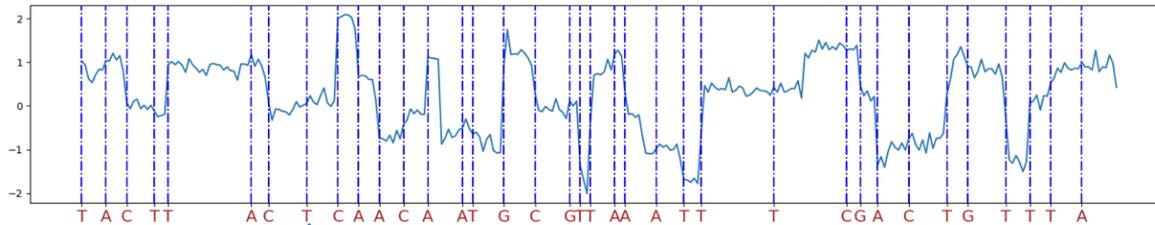
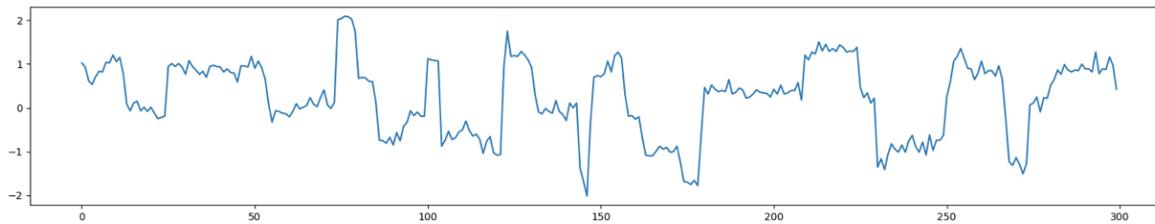
Zhang et al., PLoS Compt Biol, 2021

我々が開発したDNAシーケンスデータから変異を検出する手法の例を示します。図で示されていますが、それぞれ数式で記述される数理モデルとなっています。

研究（1）：深層学習モデル(AI)によるDNA配列推定



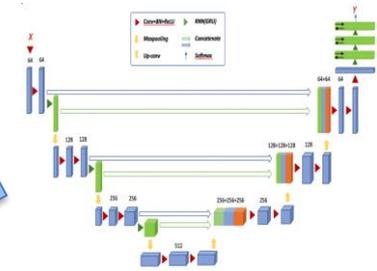
- ナノポアシーケンサーと呼ばれる、計測装置から得られるデータの解析手法の紹介をします
- 従来の装置に比べて、長いDNA断片(1000塩基以上)を計測できる反面、まだエラーが大きい問題があります
- 計測データ（電流値）の複雑なパターンからDNAの塩基を推定するのが難しいからです
- 我々は新しい深層学習モデル(URnano)を開発しました
- AIは大量のデータを学習することで、複雑な問題を解決することができます



計測データ
(電流値)

AIの予測
塩基ラベル

AIモデルURnano



AIによる高精度の予測を達成

正解
塩基ラベル Zhang et al., BMC Bioinformatics, 2019

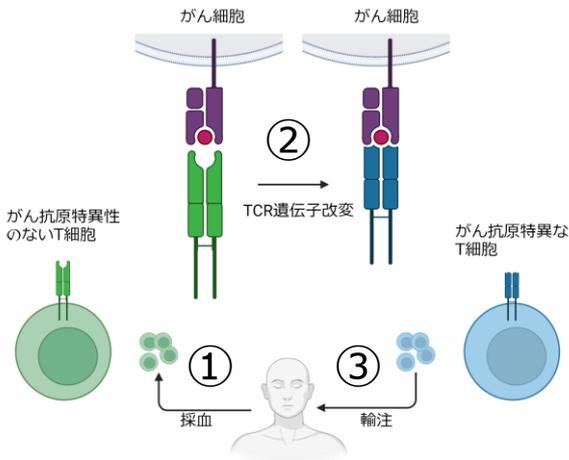


研究（2）：タンパク質結合能予測AIモデル

改変型T細胞輸注療法



結合能を高めた受容体タンパクをAIで生成したい：
ベイズ統計を応用 結果から原因を推測



Bayesian Method for TCR Design

Aim: Search for TCR structures with higher binding affinity to pMHC

受容体タンパク質

Forward Model

$$A_{TCR} = f(S_{TCR})$$

$$p(A_{TCR} | S_{TCR})$$



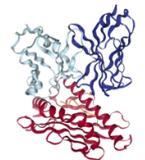
結合能予測

Binding Affinity

$$A_{TCR}$$

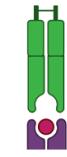


Template-based Modeling Software: Rosetta



TCR Structure

S_{TCR}



pMHC



Reverse Model

$$S_{TCR} = f^{-1}(A_{TCR})$$

Bayes' Law

$$p(S_{TCR} | A_{TCR}) \propto p(A_{TCR} | S_{TCR}) \times p(S_{TCR})$$

Generate and search TCR structures by Monte Carlo sampling



免疫細胞は受容体タンパクで、がん細胞表面の目印となるタンパクを認識して攻撃します。

AIを使って、目印となるタンパク質への結合能を高めたタンパク質をデザインする手法の開発を目指しています。

研究（２）：タンパク質結合能予測AIモデル

前のページで述べた結合能を高めた受容体タンパクの設計に向けて、まずタンパク同士の結合能力を、アミノ酸配列情報と立体構造の情報から正確に予測するAIの開発に成功しています。

A multimodal framework for protein-protein binding affinity prediction

タンパク質配列情報 ΔG of a wild type complex and that of a mutant complex

Wild type sequence Chain A: HPETLVKVKDAED
Chain B: AGVMTGAKFTQIQ

mutant sequence Chain A: HPETLVAVKDAED
Chain B: AGVMTGAKFTQIQ

Siamese-ESM network

Sequence feature

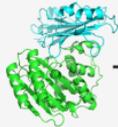
結合能予測

Gradient boosting tree

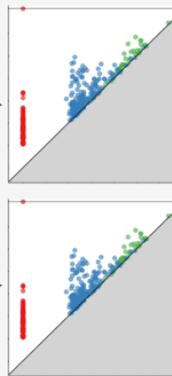
$\Delta\Delta G$

タンパク質立体構造情報

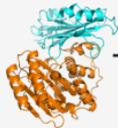
Wild type structure



Topological data analysis



mutant structure



Topological data analysis

Convolutional neural network

Structure feature

研究（3）：情報管理・解析システムの開発

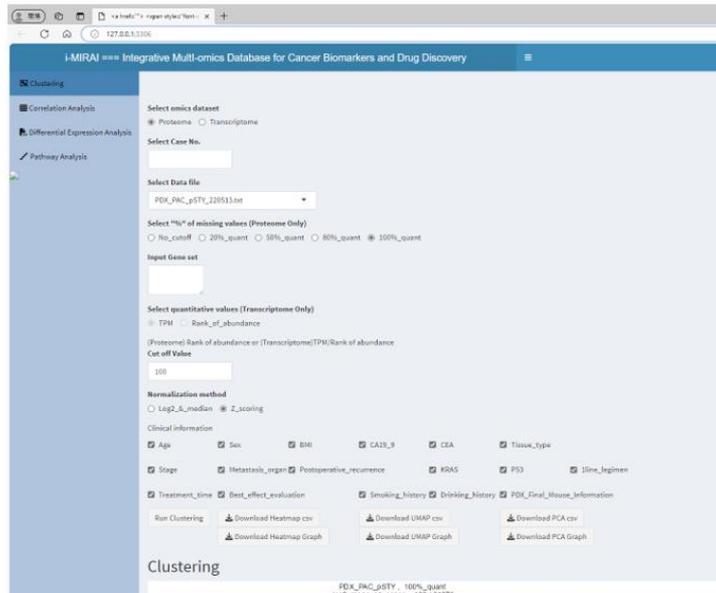


重点プロジェクトで産生される大量かつ複数種類のデータを、特別なプログラミングなどを必要とせずに統合的に解析できるシステムの開発を行っています。

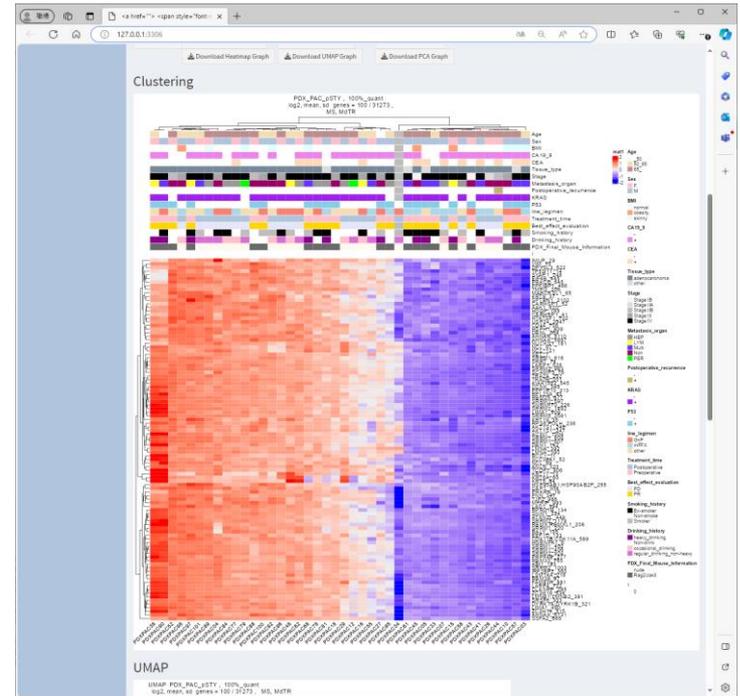
データ管理・共有・統合・再活用
Information management system



INPUT 項目 ←

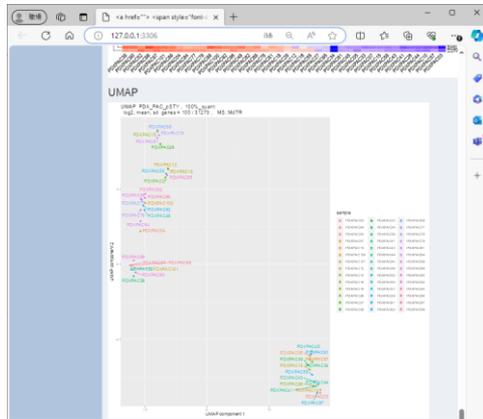


クラスタリング/Heatmap

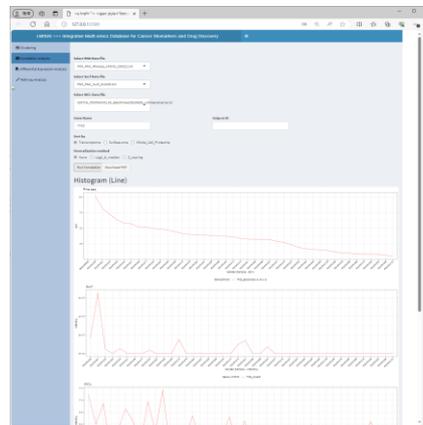


実際の研究現場では複雑なデータに対して、条件を変えながら複数種の解析を実施し、得られた多角的な情報を基に仮説を立て、次のステップへ進むことを繰り返します。このシステムはそのプロセスを加速します。

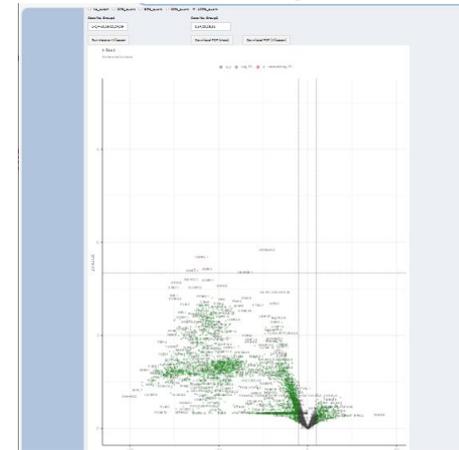
次元削減 (UMAP/PCA)



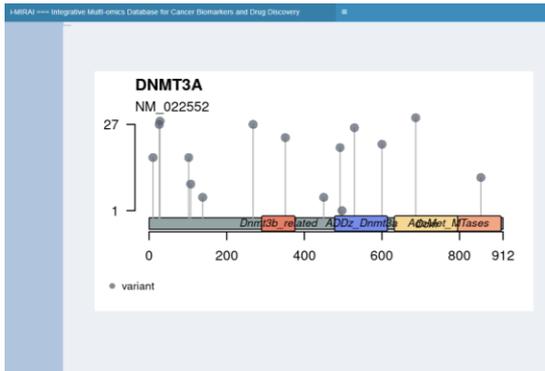
相関解析



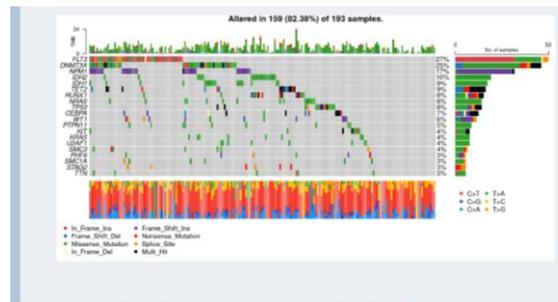
発現差解析



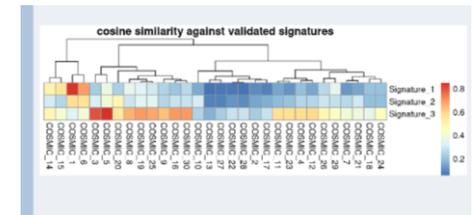
Lollipop plot



Oncoplot



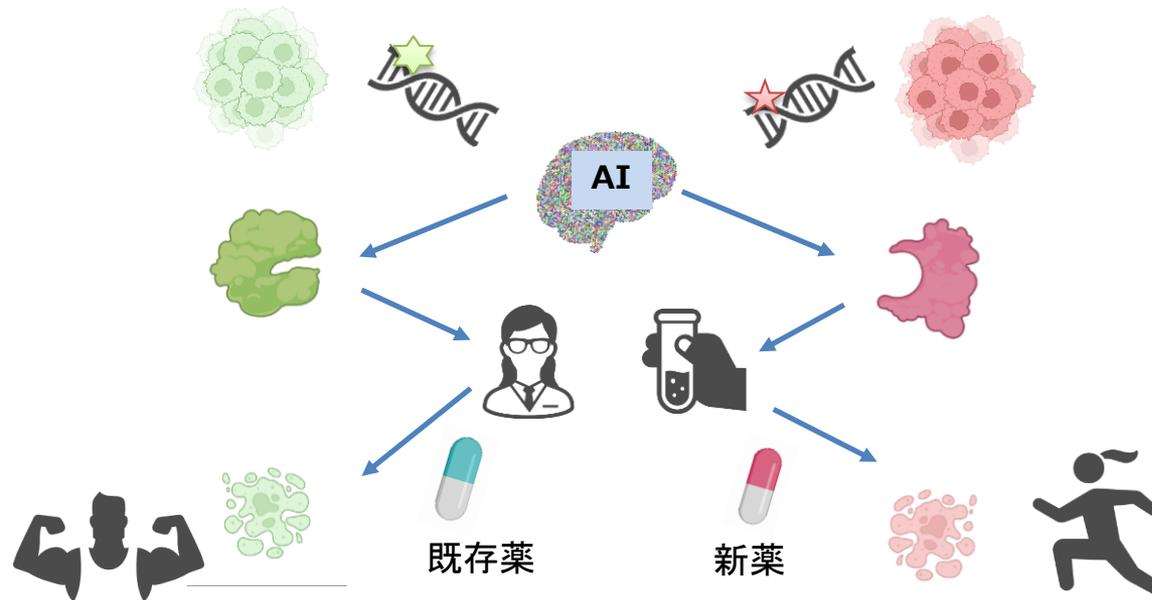
Mutational signature



関連した話題：タンパク質の立体構造がAIで予測できると

実験的に決定困難なタンパク立体構造を
これまでにない速度と労力で推測することができる

がんのメカニズムの究明や薬の開発が加速度的に進むことが期待



タンパク質の立体構造と性質は密接に関わっています。しかし多くのタンパク質の立体構造を実験的に決めることが難しいことが課題でした。

その課題に対して近年、タンパク質のアミノ酸配列（1次元情報）から立体構造（3次元情報）を予測するAIの開発が、世界中で急速に進んでいます。

高校生向け基礎実験体験講座

2022 Summer Seminar 高校生向け基礎実験体験講座

AIでがんを観る!

～ スパコンで生命医学研究を体験しよう～

開催日時
2022.8.5(金)
9:00-17:00

募集人数
終日参加出来る高校生、計15名

開催場所
愛知県がんセンター

お申し込み方法
事前申込みが必要です。
申込期間：7月27日(月)から7月8日(金) (別印有効)まで
はがき(別印有効)、FAXでお申し込み下さい。
その他、「高校生体験講座希望」、氏名、心のがな、郵便番号、住所、電話番号、生年月日、性別、学校名、学年を明記してください。

☆参加者へのお知らせ
事前申込者の中から、抽選で参加者を決定いたします。
参加の可否は7月18日(月)以降に郵送いたします。

タンパク質の分子構造を可視化したもの

実習概要

本実習ではAlphaFoldという人工知能プログラムを利用して、タンパク質のアミノ酸配列から立体構造を予測します。

これまでこのべ**244人**がこの体験講座に参加されました

参加者の声

とても楽しかったです。なかなか体験できないことを体験できました。

見たことない構造がいっぱいあっておもしろかったです!

がんや色々な病気に関係を持つことができました。

研究者の人たちの考えや話が聞けてよかったです。

愛知県がんセンターでは毎年8月に高校生を対象とした体験講座を実施しています。
一昨年は、将来の研究者の卵である高校生達に最先端のタンパク質立体構造予測AI (AlphaFold*) を体験してもらいました。

*: Jumper et al., Nature, 2021

高校生がいち早く最新のAIを体験!



講義の様子1 : ⇐

AlphaFold のアルゴリズムについての講義⇐



実習の様子1 : ⇐

スパコンへのログイン⇐



実習の様子3 : ⇐

AlphaFold から予測された EGFR タンパク立体構造の可視化。
野生型と変異型の構造を重ね合わせて比較。⇐

おわりに

次世代シーケンサーによりがん細胞から取得される大規模計測データ（がんビッグデータ）とその多様化について紹介しました

個々のがんの性質を知り治療戦略の立案につなげるには、それらのデータから有用な情報を抽出するための手法およびシステムの開発が重要です

ここではシステム解析学分野で開発を進めている解析手法および解析システムの一端を紹介しました

医療AIの開発と発展が進むなか、当分野では愛知県がんセンターの病院と研究所の方々と協力して、データ科学の力で医療に貢献できるよう研究を進めてまいります