

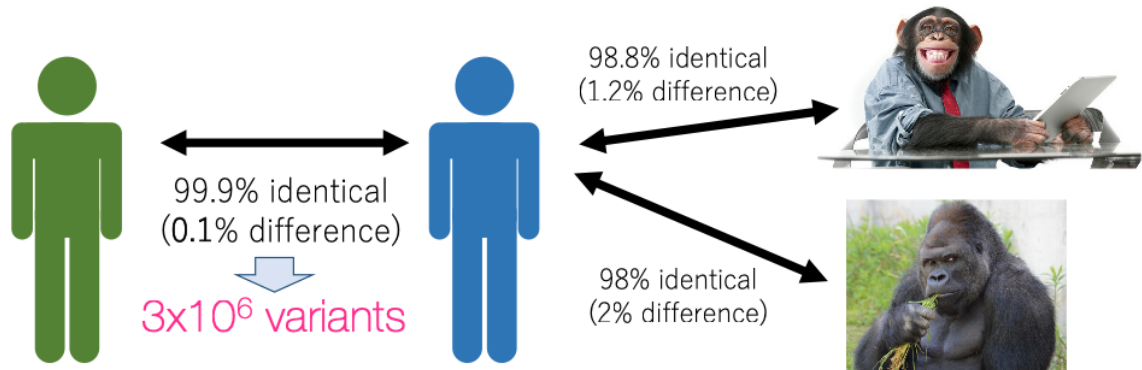
# ゲノムバリエントと表現型との 関係を分子構造情報でつなぐ データサイエンス

土方 敦司

東京薬科大学生命科学部 生命医科学科  
ゲノム情報医科学研究室

221209 東北大学未踏データアナリティクスセンター  
第2回UDACセミナー

表現型の違い(個人差・多様性)の情報は  
ゲノムにどのようにエンコードされているか？



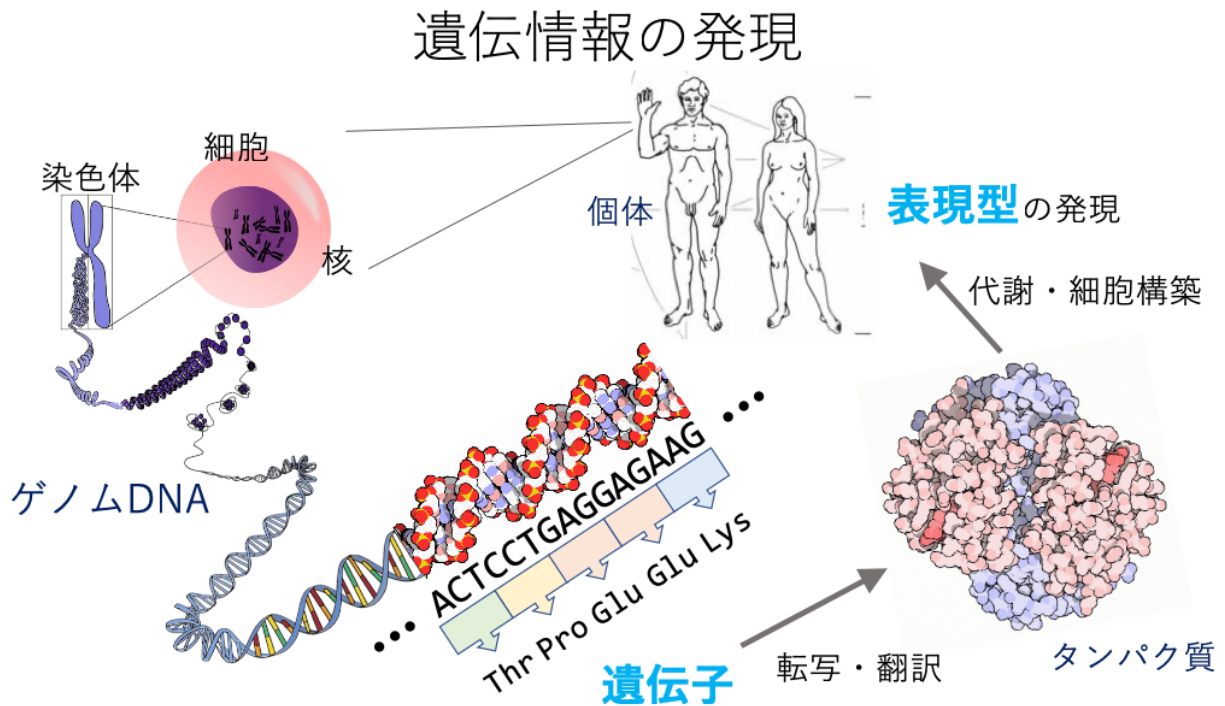
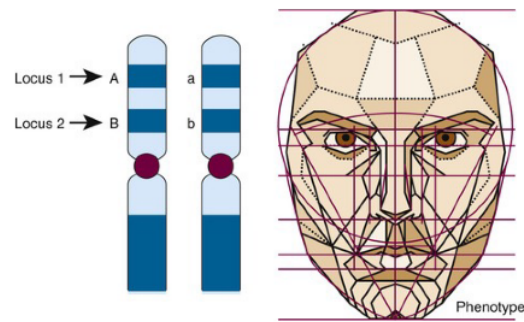
ゲノムは「生命の設計図」であるので、ゲノムの類似度は生物の近縁関係を反映  
同一ヒト集団内のゲノム上の差異(バリエント)が個性を生み出す一因と考えられる

ただし、全てのバリエントが表現型と関係しているわけではない

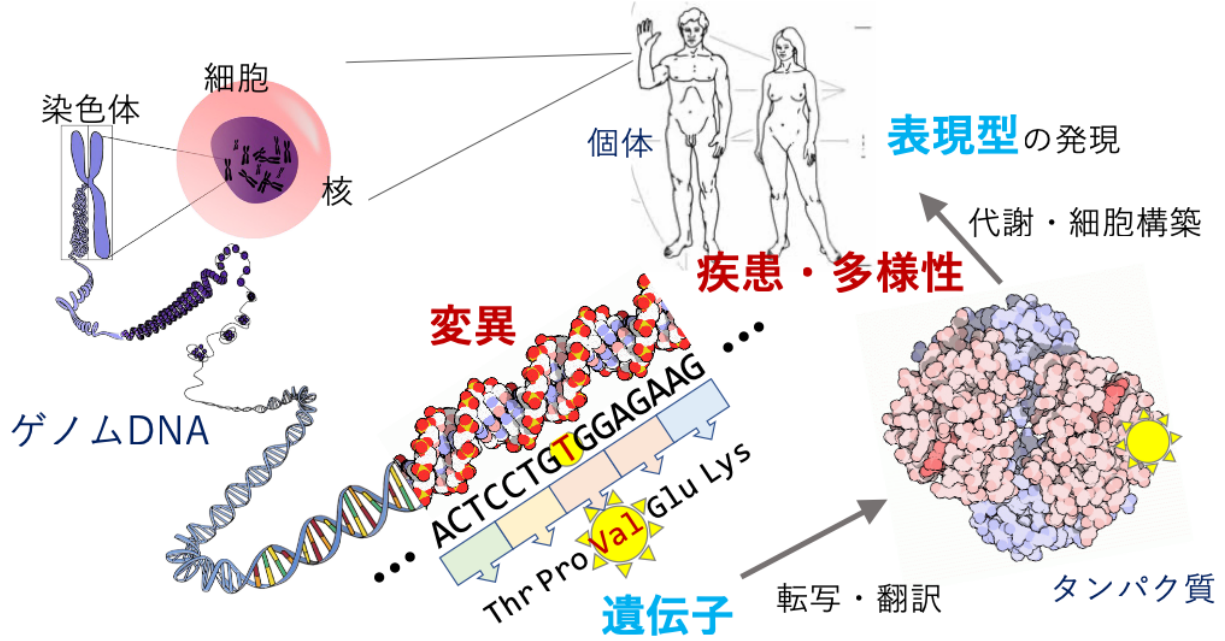
# 表現型(Phenotype)

遺伝型の発現。形態学的、臨床的、細胞学的もしくは生化学的な**形質**として臨床的に観察できるか、または血液や組織検査でのみ検出できる:

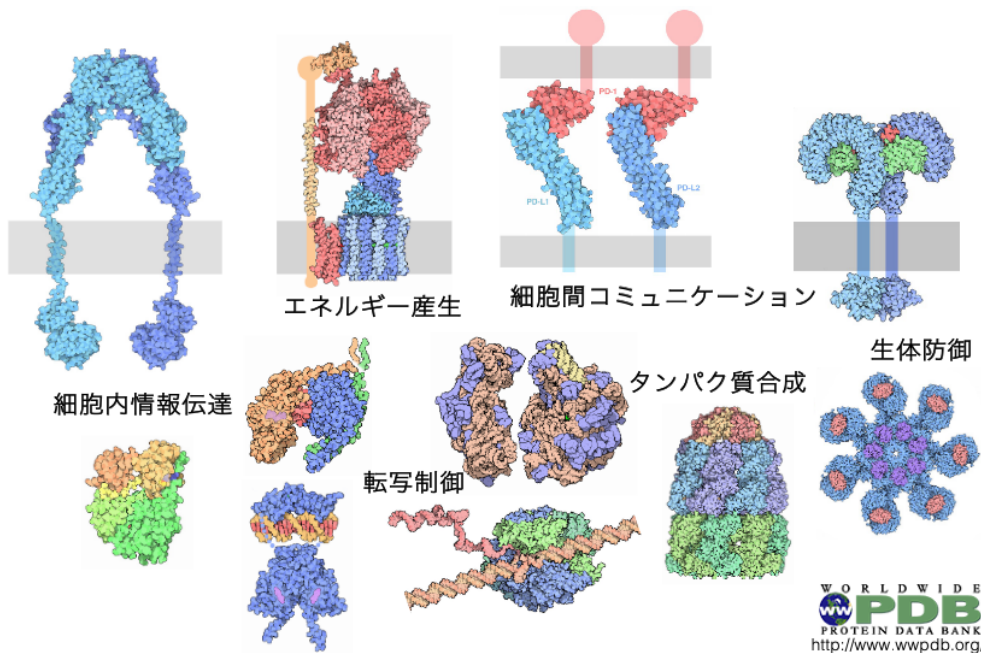
例：疾患、身長、目の色、顔の形、血液型など



# 遺伝情報の発現

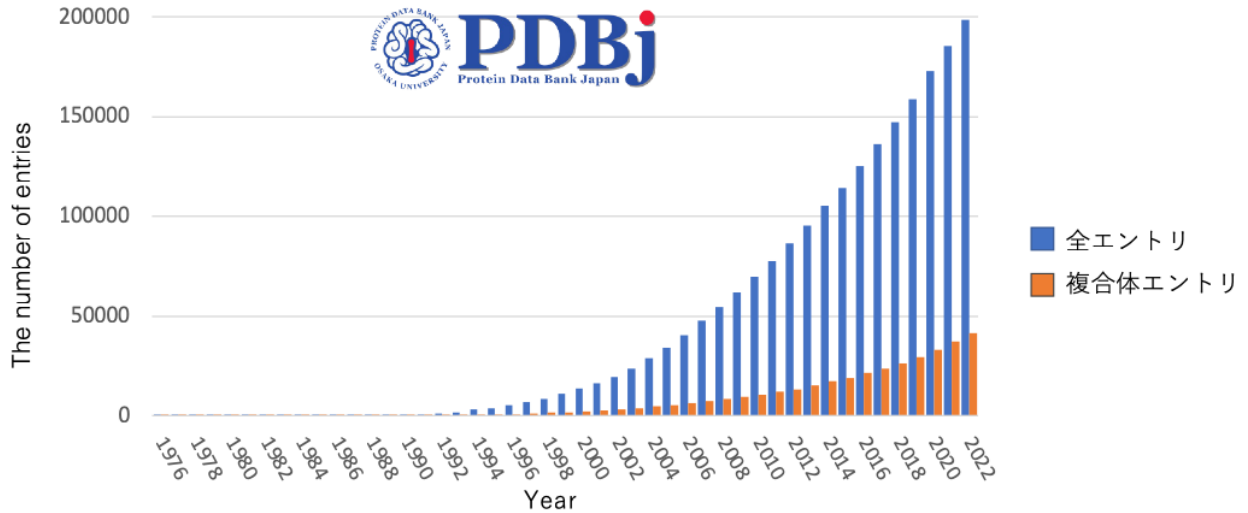


タンパク質は生体内で“超分子複合体”として機能する

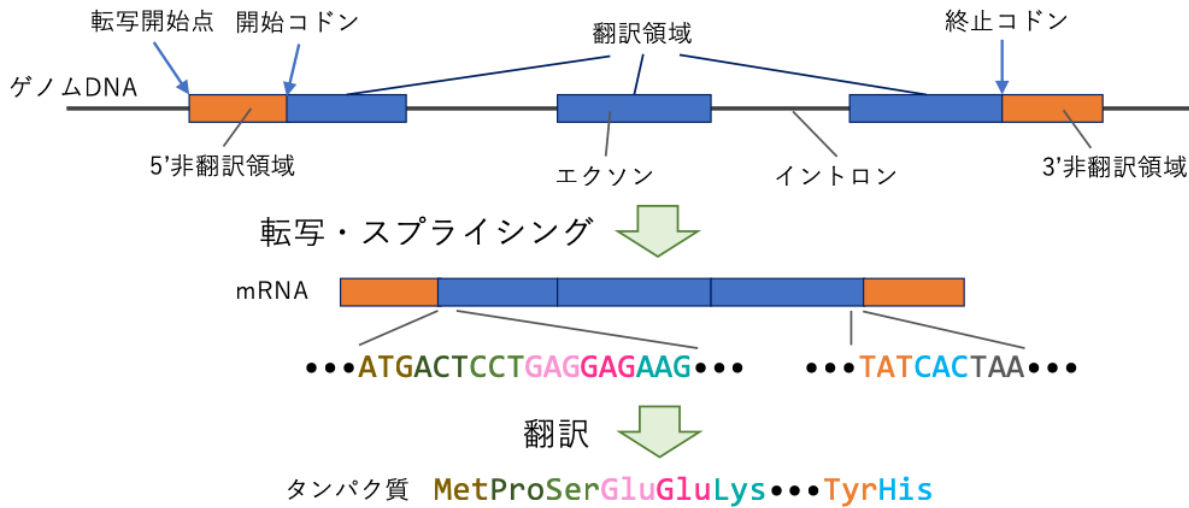


# 分子構造情報の蓄積状況

- Protein Data Bankに登録されている構造データは20万に迫る(2022年11月時点)
- 複合体の構造データも増加してはいるものの、完全には程遠い



# 真核生物の遺伝子構造





# ゲノムバリアントの種類

一塩基置換

ATAC**GT**GCTA  
ATAC**CG**GCTA

多塩基置換

ATAC**GT**GCTA  
ATAC**CCG**GCTA

挿入

ATACG--TGCTA  
ATACG**CG**TGCTA

欠失

ATACG**T**GCTA  
ATACG-**G**CTA

重複

ATACG-----**TG**CTA  
ATACG**TGCTGCTG**CTA

## ゲノムバリアントのタイプとタンパク質への影響



DNA ...ACTCCTGAGGAGAAA...

タンパク質 ...ProSerGluGluLys...

ミスセンス変異

一塩基置換

...ACTCCTG**T**GAGAAA...  
...ProSer**Val**GluLys...

ナンセンス変異

...ACTCCTGTGGAG**TAA**...  
...ProSerGluGlu**Stop**...

フレームシフト変異

挿入(欠失)

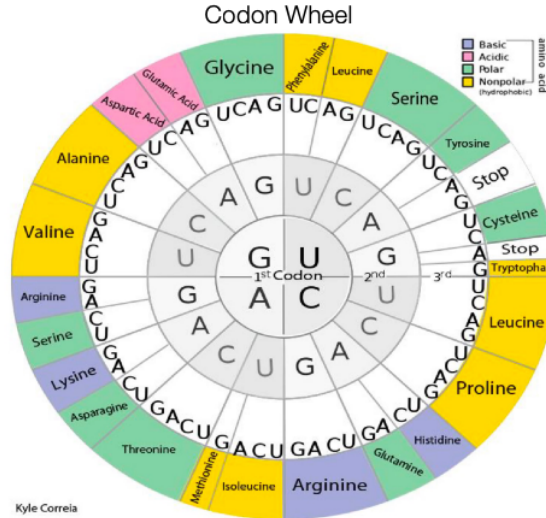
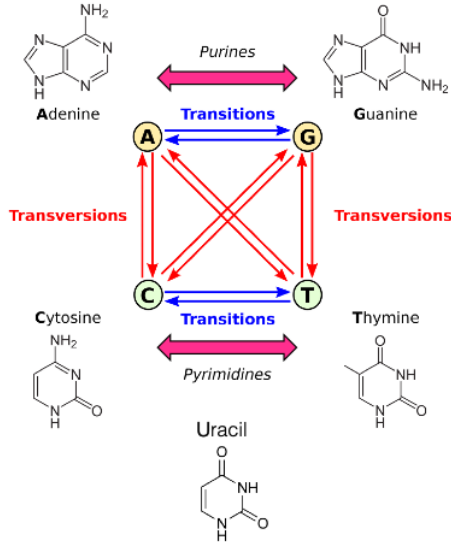
...ACTCCTG**AC**GAGAAA...  
...ProSer**AspGly**Glu...

インフレーム挿入(欠失)

...ACTCCTGAG**TTT**GAGAAA...  
...ProSerGlu**Phe**GluLys...

# 塩基置換(ゲノム)とアミノ酸置換(タンパク質)との関係

- 塩基置換はプリン同士、ピリミジン同士が起こりやすい性質
- アミノ酸とコドンの対応は塩基置換(エラー)に対して頑健



# 遺伝疾患の遺伝形式と発現形式について

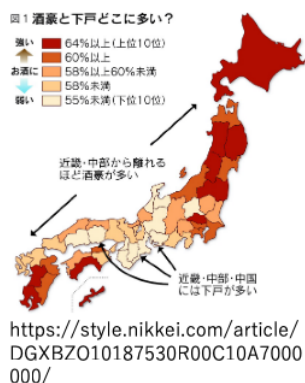
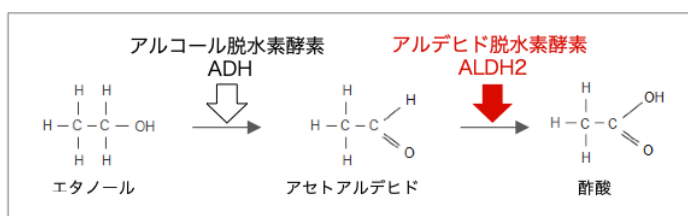
- 顕性(優性)遺伝の疾患発現形式は、比較的複雑でありメカニズムの解明があまり進んでいない
- 単一遺伝子の変異であっても、複数の遺伝形式や発現形式をとる場合がありうる

	遺伝形式	発現形式	遺伝子型
Recessive 潜性(劣性)	AR (Autosomal Recessive) 常染色体劣性	LF (Loss-of-Function) 機能欠損	 Homozygous(ホモ)
	XLR (X-linked Recessive) 伴性(X連鎖性)劣性		
Dominant 顕性(優性)	AD (Autosomal Dominant) 常染色体優性	HI (Haploinsufficiency) ハプロ不全	 Heterozygous ヘテロ
		DN (Dominant Negative) 優性阻害	
	GF (Gain-of-Function) 機能獲得		
XLD (X-linked Dominant) 伴性(X連鎖性)優性			
非メンデル性遺伝		ミトコンドリア遺伝 (母系遺伝)	

## 単一遺伝子のバリエーションによって表現型が決まる例

表現型 (形質・疾患)	原因遺伝子	発現形式
お酒が飲めるか (アルコール感受性)	<i>ALDH2</i>	顕性(AD)
耳垢のタイプ(ベタベタ or カサカサ)	<i>ABCC11</i>	顕性(AD)
ワルファリン(抗血液凝固剤)の応答性	<i>CYP2C9</i>	顕性(AD)
ハンチントン舞蹈病	<i>HTT</i>	顕性(AD)
フェニルケトン尿症	<i>PAH</i>	潜性(AR)
鎌状赤血球症	<i>HBB</i>	潜性(AR)

## ALDH2：お酒の強さとの関係



**Glu504**  
 ALDH2\*1 ...ACT-**G**AA-GTG... 正常型アレル  
 ALDH2\*2 ...ACT-**A**AA-GTG... 低活性型アレル  
**Lys504**

rs671 (p.Glu504Lys)

遺伝子型	お酒に強い		お酒に弱い	
	G	G	G	A
日本人の割合 <sup>a)</sup>	67 %		30 %	3 %
欧州人の割合 <sup>b)</sup>	100 %		0 %	0 %

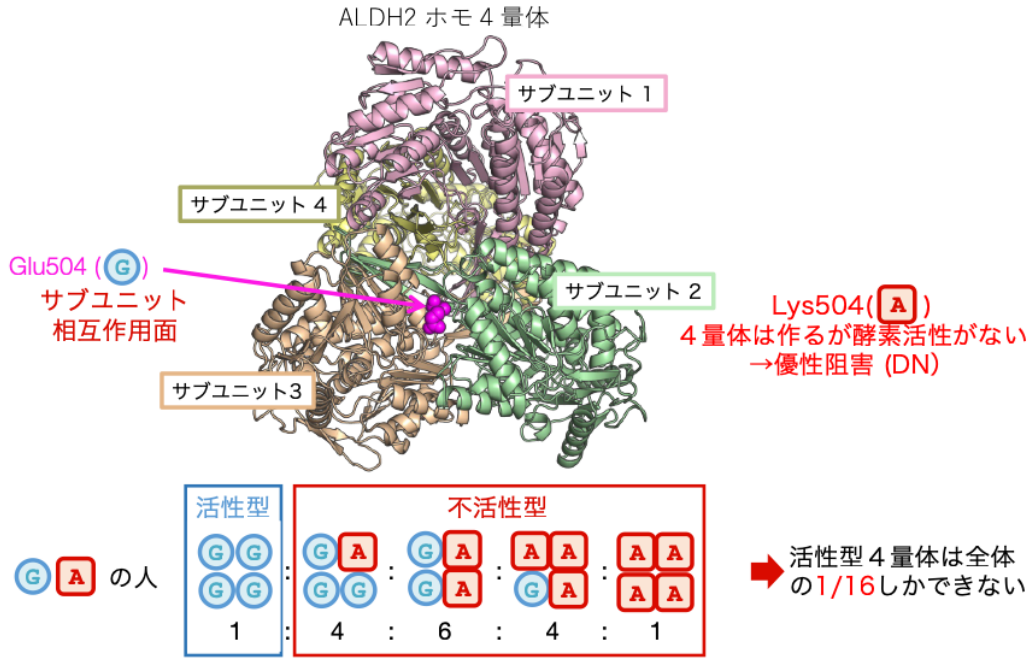
ALDH2活性比

$$\frac{\text{G} \text{ A}}{\text{G} \text{ G}} \rightarrow \frac{1}{16}$$

なぜ1/2でないか？

a) <https://ijgvd.megabank.tohoku.ac.jp>  
 b) <http://www.internationalgenome.org>

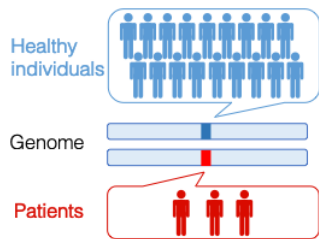
# ALDH2はホモ 4 量体を形成して機能する



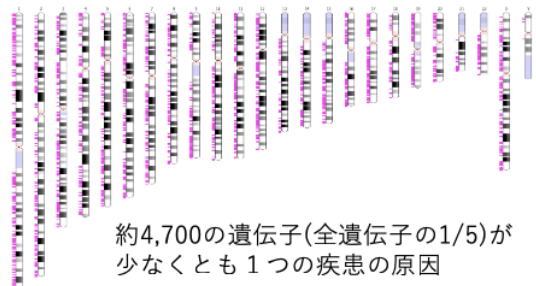
# NGSの登場による疾患変異同定の加速



## 患者群と健常者群のWES/WGS

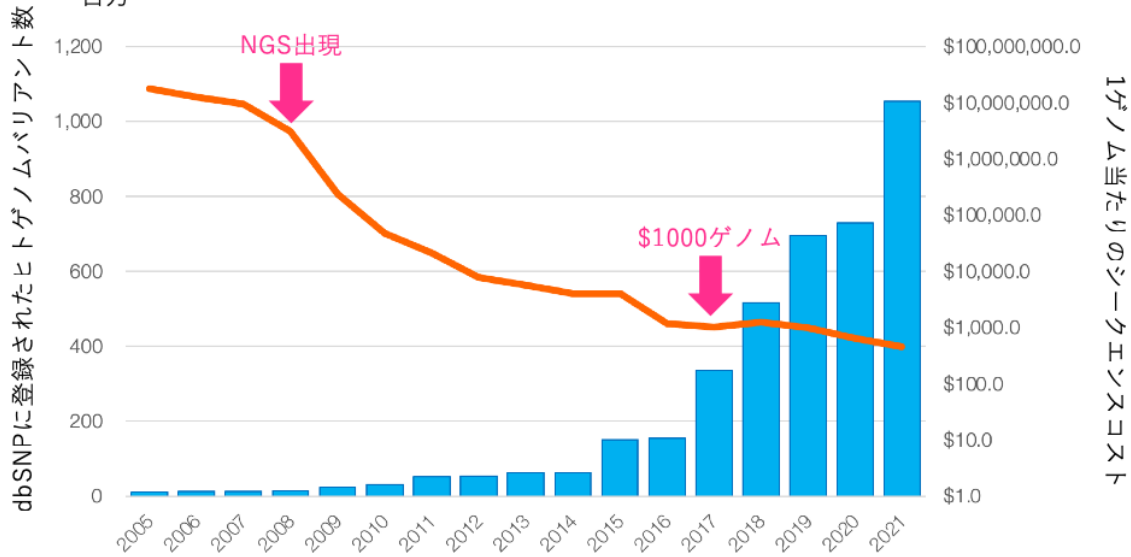


## 疾患原因遺伝子のゲノムマップ

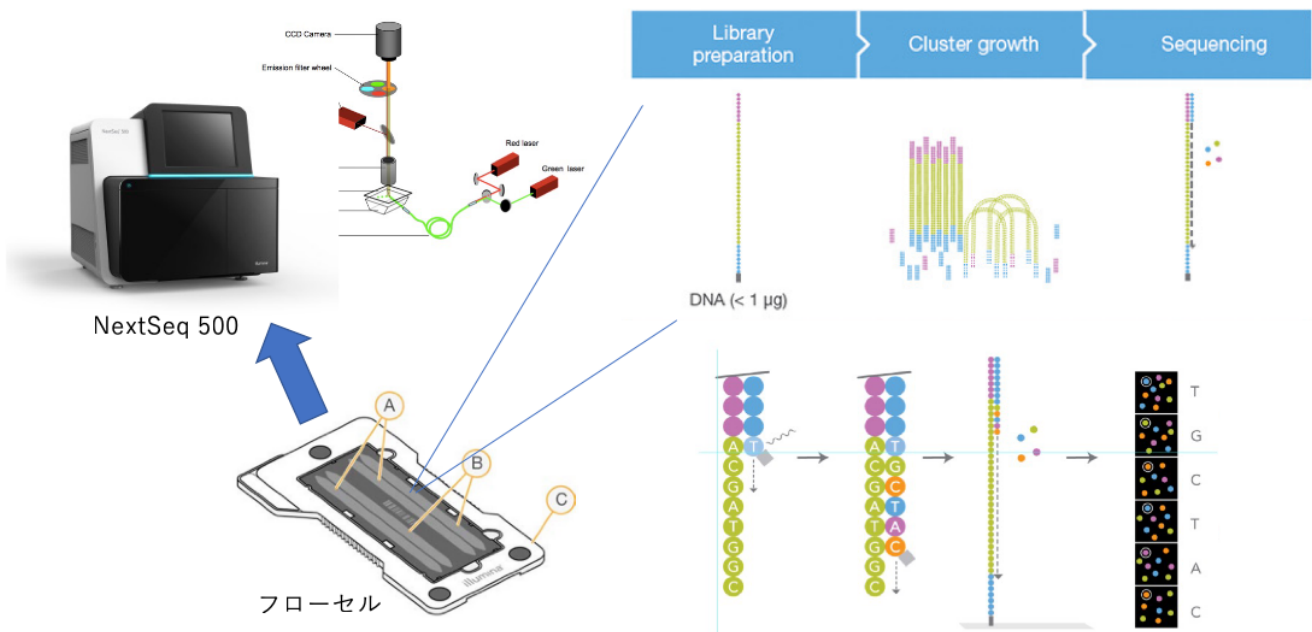


# ゲノムバリエーションデータの爆発的増加

シーケンスコストの減少に伴いヒトゲノムバリエーションデータは増加の一途を辿っている  
百万



# NGS(次世代シーケンサー)の原理

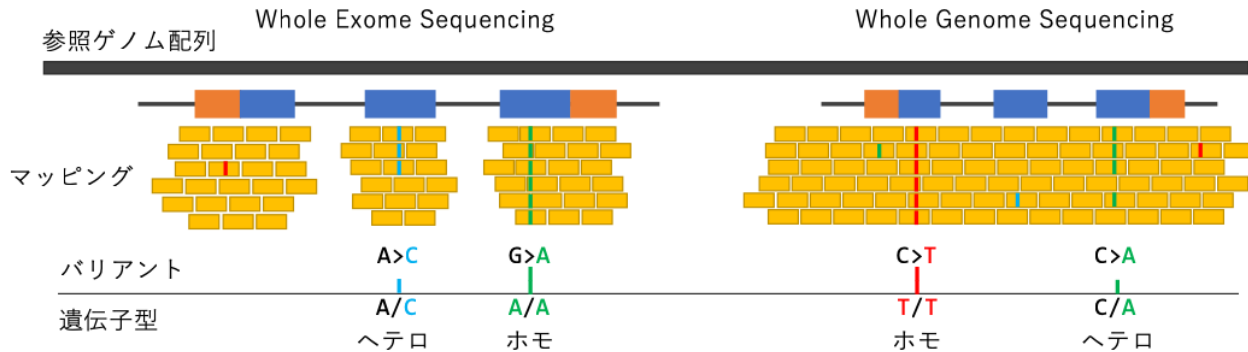




# NGSによるゲノムバリアントの同定

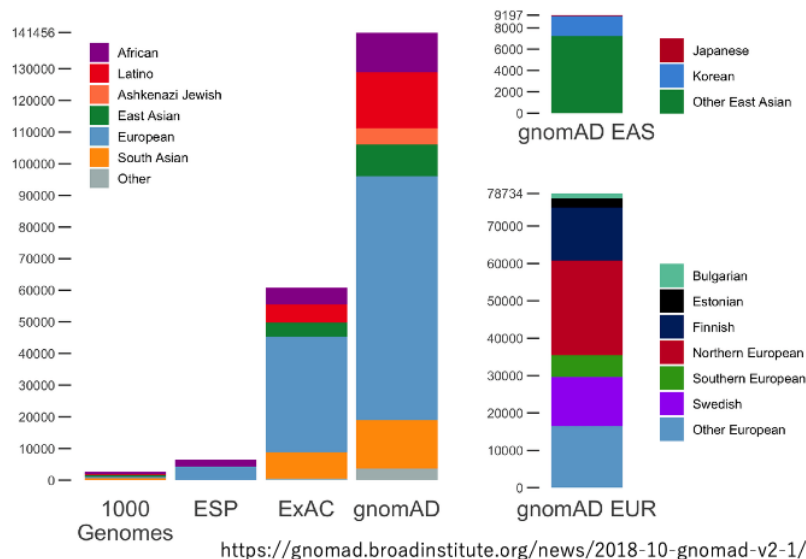
- NGSから出力される短いDNA配列(100-200nt)を、参照ゲノム配列に「マッピング」する
- 参照配列と比べて「変化」しているところを読み取りエラーと区別してバリアントとしてコールする
- 全ゲノム配列を読みとる場合と全エクソン(Exome)を読みとる 2種類がある
- ヒトゲノムは2倍体(父親由来のアレルと母親由来のアレルがある)ため、バリアントはホモ(両アレルともバリアント)とヘテロ(片アレルのみバリアント)で検出される

■ NGSから出力されるDNA配列(リード)



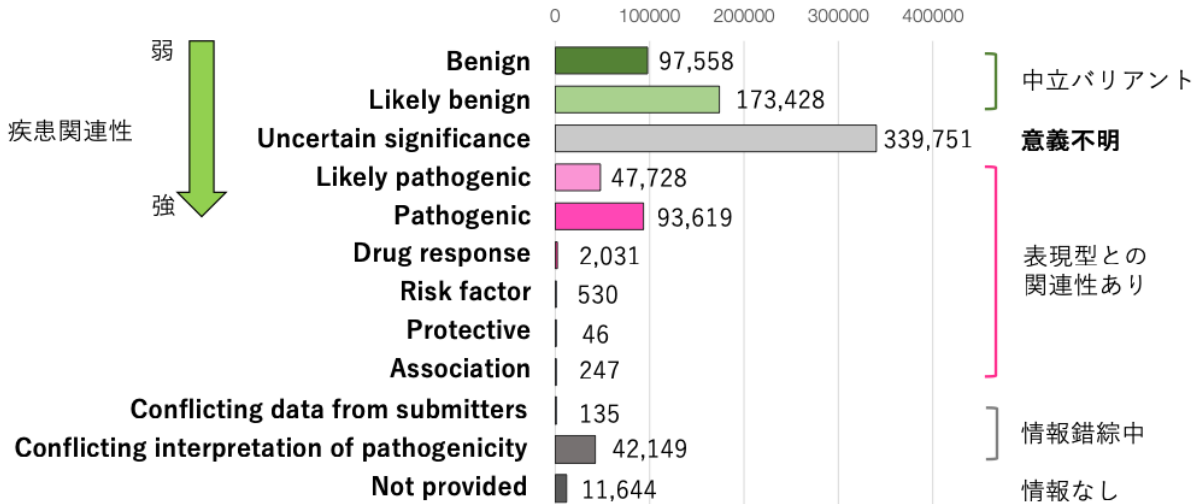
## パーソナルゲノムデータの蓄積が進んでいる

- 健常人に見られる各バリアントの頻度情報が公開(2021年時点で約7億バリアント: 4塩基に1箇所)
- バリアントの頻度情報は、疾患変異とそうでない変異を区別する上で重要な役割を持つ
- 一方で、表現型との関連がよくわからない低頻度バリアント(レアバリアント)も多い



# ClinVar

- ヒトバリエントと(疾患)表現型の関連性データベース <https://www.ncbi.nlm.nih.gov/clinvar/>
- 研究者が基準に従って臨床的有意性(Clinical significance)をつけている
- 意義不明なバリエント(VUS)が突出して多い



## NGS(配列)ではわからないこと

ゲノムバリエントと疾患との“関連性”はわかったが“メカニズム”は？



変異による疾患のメカニズムを理解することは、

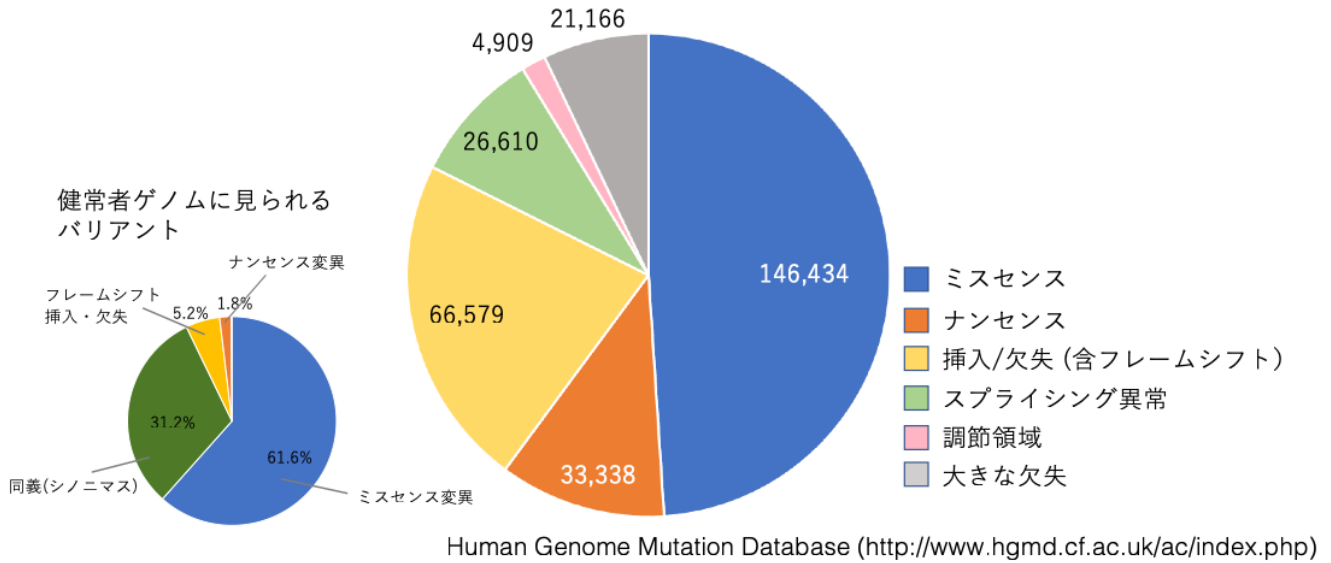
正確な診断や効果的な治療法の確立や  
オーファンドラッグ(希少疾患用医薬品)の開発に重要



遺伝子がコードするタンパク質の機能は？  
そのアミノ酸残基の役割と変異の影響は？

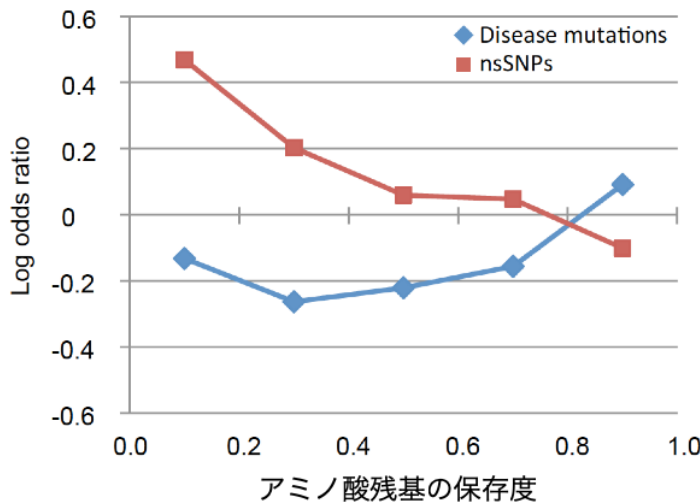
# 疾患原因バリエーションの種類

- 既知の疾患原因バリエーションの約半数はミスセンスバリエーションが占める。
- 一方で、疾患と関係のない中立なミスセンスバリエーションも多い。どのように両者を見分ければよいか



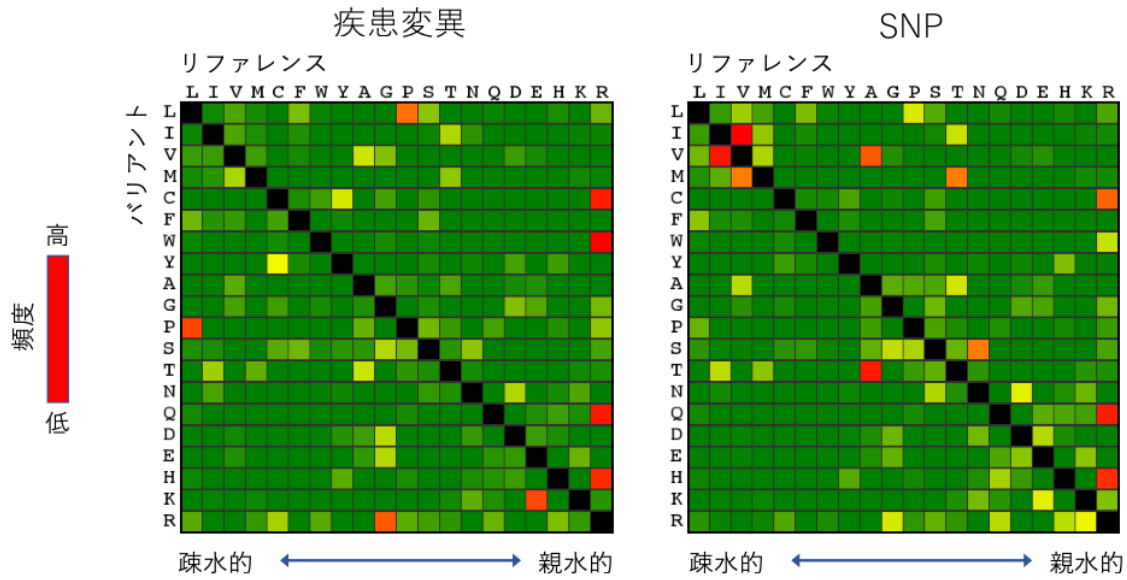
## 疾患変異は進化的に保存されたアミノ酸残基に起きやすい

タンパク質の機能や構造にとって重要なアミノ酸残基は生物種間で保存される傾向にあるため、そのようなアミノ酸に変異が入ると機能を損なってしまう可能性が高い

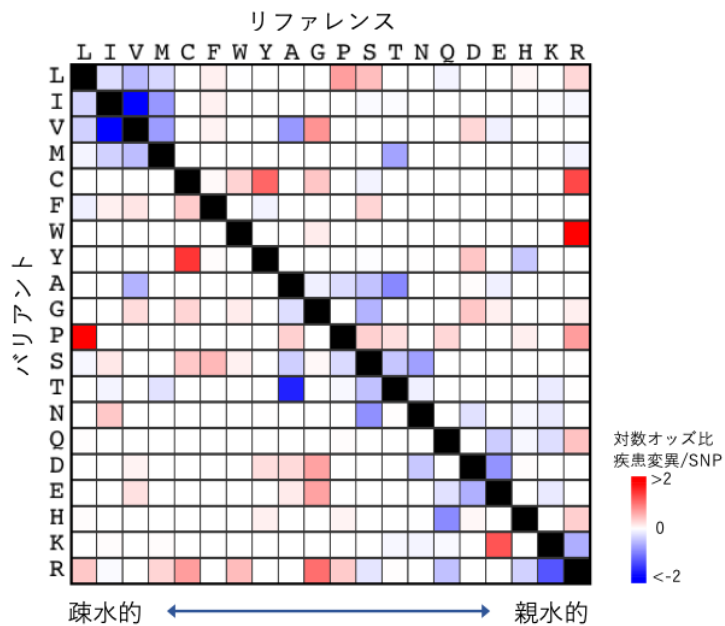


Hijikata et al. (2010)

疾患変異はどのようなアミノ酸置換パターンが多いのか？

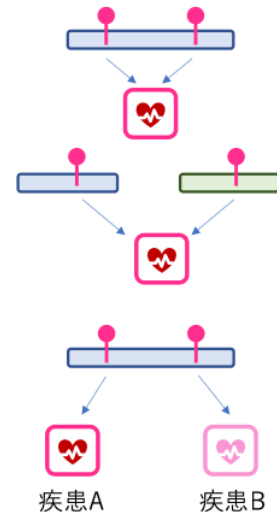


疾患変異はどのようなアミノ酸置換パターンが多いのか？



# 遺伝疾患における異質性

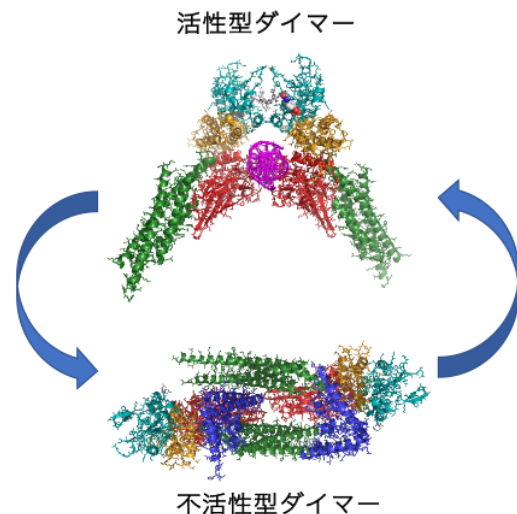
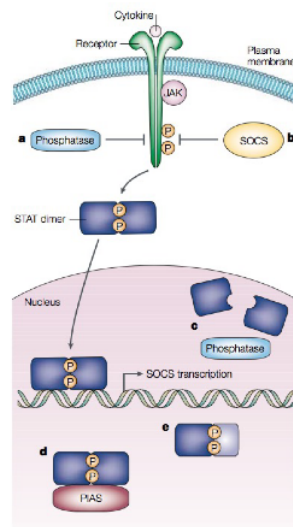
- ◆アレル異質性(Allelic heterogeneity)  
同一遺伝子の異なる変異が同一の疾患の原因となる
- ◆座位異質性(Locus heterogeneity)  
異なる遺伝子の変異が同一の疾患の原因となる
- ◆表現型異質性(Phenotypic heterogeneity)  
同一遺伝子の異なる変異が異なる疾患の原因となる



STAT1

## STAT1 (Signal Transducer and Activator Transcription factor 1)

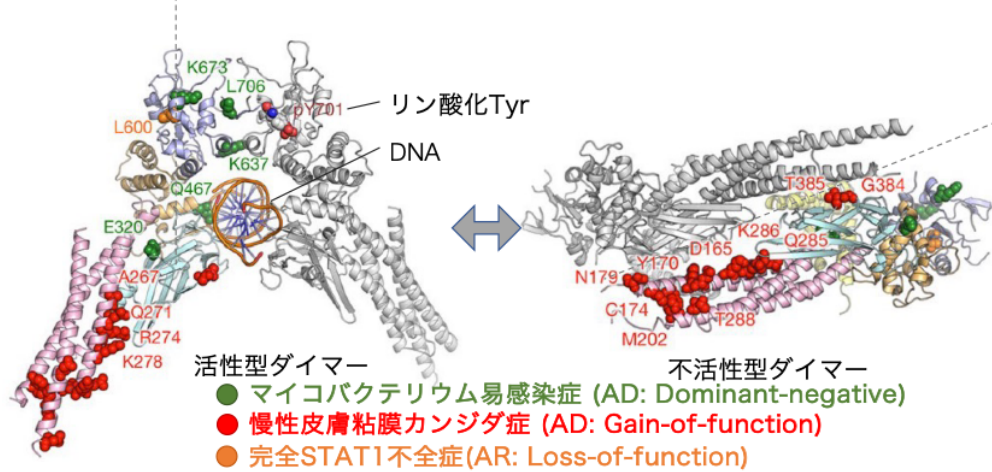
- サイトカインシグナルを下流に伝えるシグナル因子かつ転写因子
- リン酸化(活性化)されると、核に移行しDNAに結合し、遺伝子発現を制御





## STAT1は変異の違いによって異なる疾患病態をとる

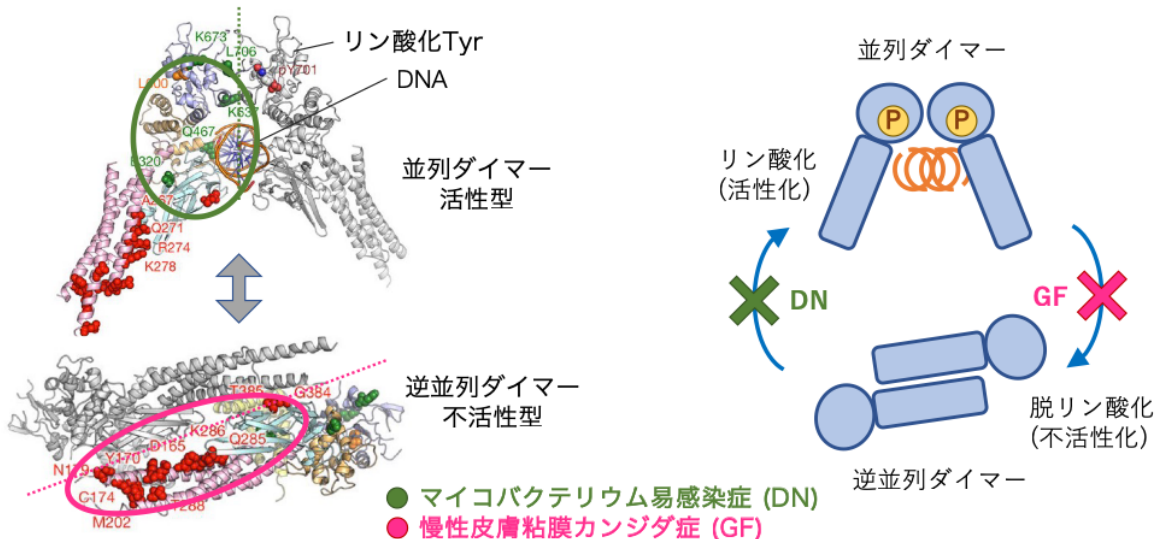
疾患名	遺伝形式	症状
完全STAT1欠損症	AR	重篤な細菌・ウイルス感染症
機能欠損型STAT1不全症	AD	結核菌への易感染性
機能獲得型STAT1不全症	AD	慢性皮膚粘膜カンジダ症・細菌感染無し



Hijikata et al. Sci. Rep. 7, 8541 (2017).

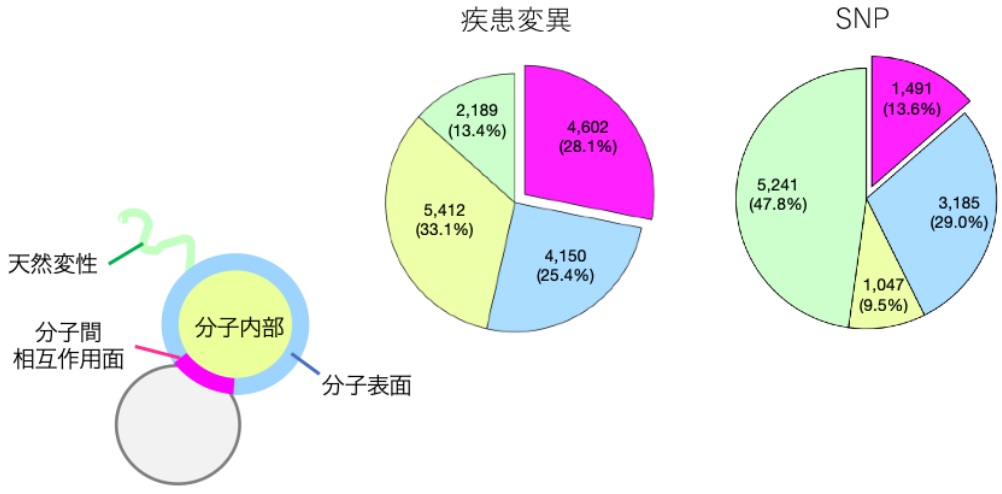
## STAT1における疾患変異の超分子複合体構造との関係

- STAT1は変異の違いによって異なる疾患発現形式をとる。
- 疾患発現形式の違いは超分子構造と密接な関係があった。



# 疾患変異はタンパク質立体構造のどの部分に多いか？

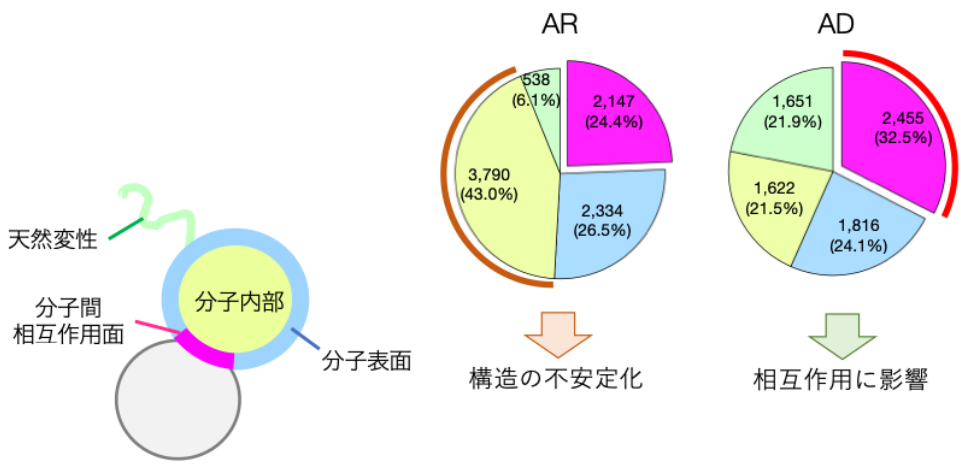
- 疾患変異はSNPに比べて、タンパク質の内部や分子間相互作用面などに有意に多く分布



Hijikata et al., Sci. Rep. (2017)

# 疾患変異を潜性(AR)と顕性(AD)に分けてみると

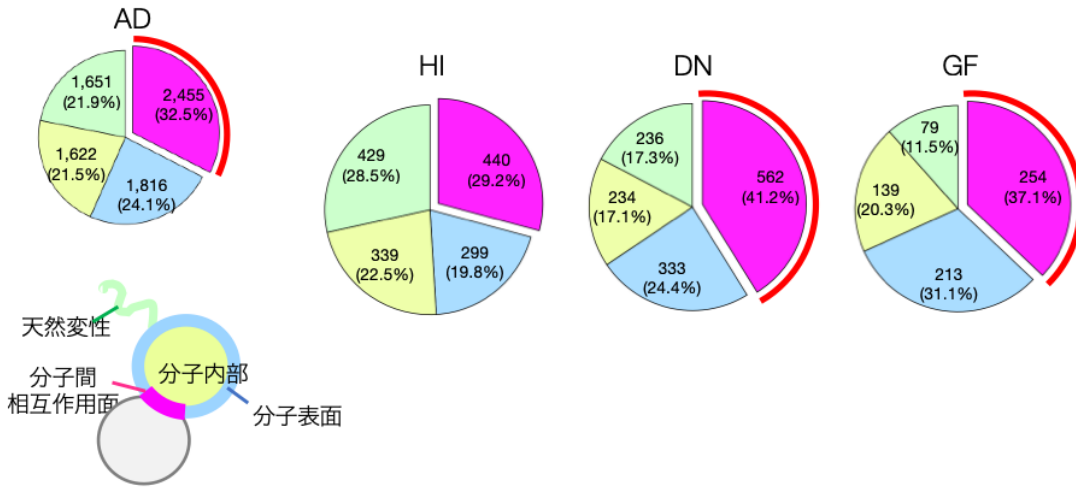
- ARの疾患変異はタンパク質内部、ADの疾患変異は分子間相互作用部位に多い



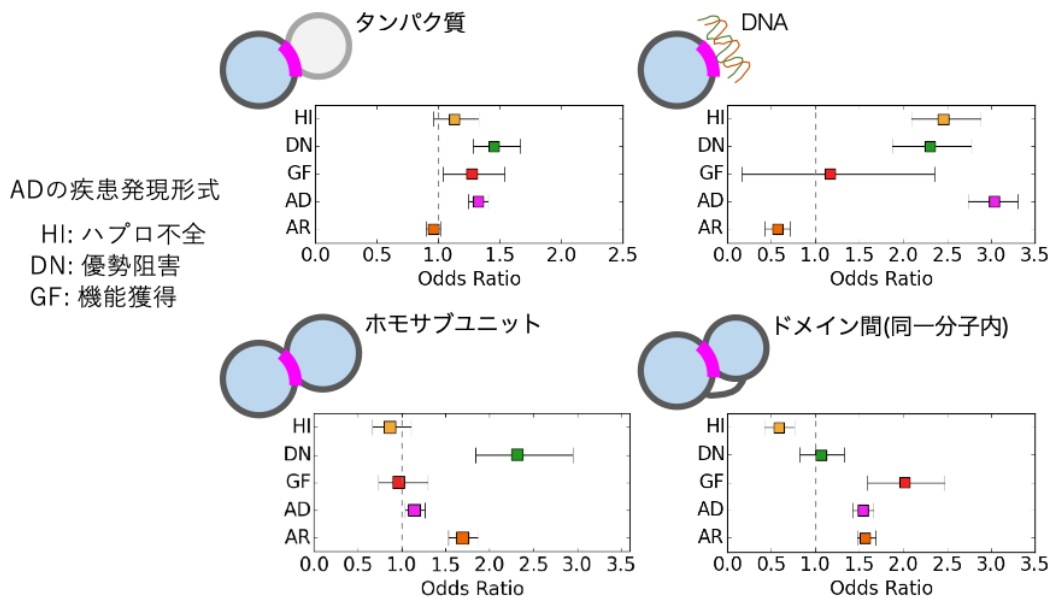
Hijikata et al., Sci. Rep. (2017)

## さらにADをメカニズムごとに分けて見ると

- 優性阻害(DN)と機能獲得(GF)ではハプロ不全(HI)に比べて相互作用部位への変異が多い傾向があった

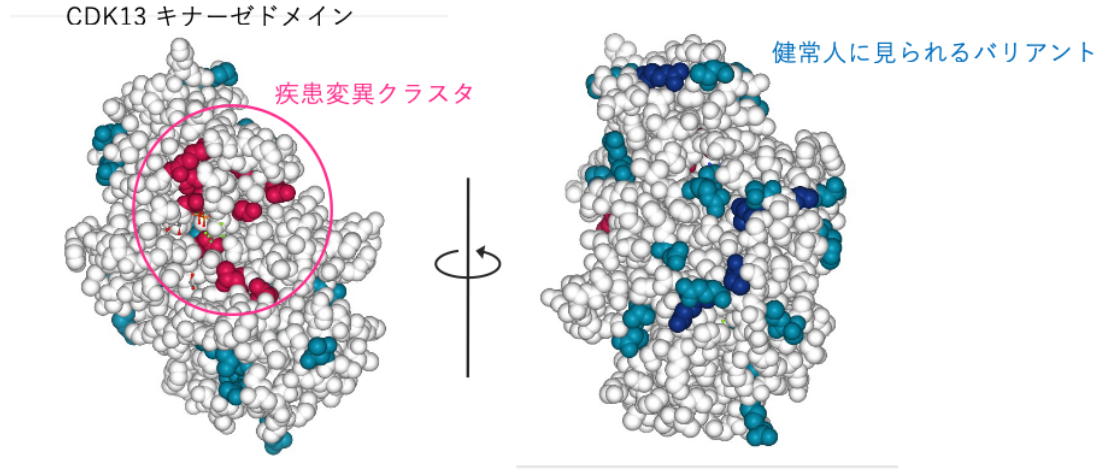


## 疾患発現メカニズムと分子相互作用タイプの関係



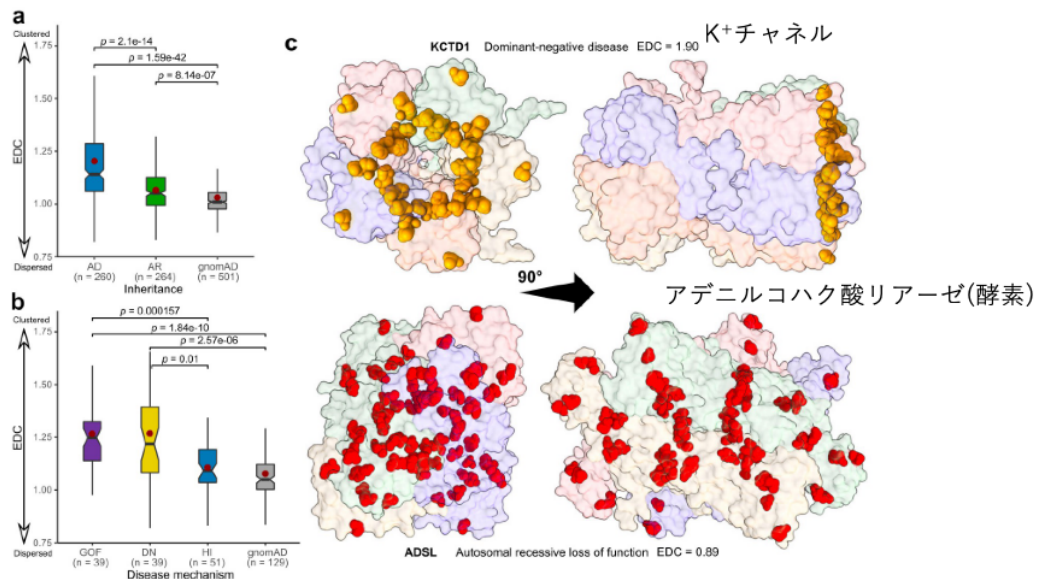
# 疾患変異は空間的にクラスタを形成していることがある

- 疾患変異はタンパク質の構造、機能に重要なアミノ酸残基に起きることが多い。
- 機能部位は、配列上離れたアミノ酸残基であっても空間的には近い位置に存在することがわかる。
- 疾患とは関係のない中立なバリエーションは、同じドメイン内であっても機能に重要でないところにある。

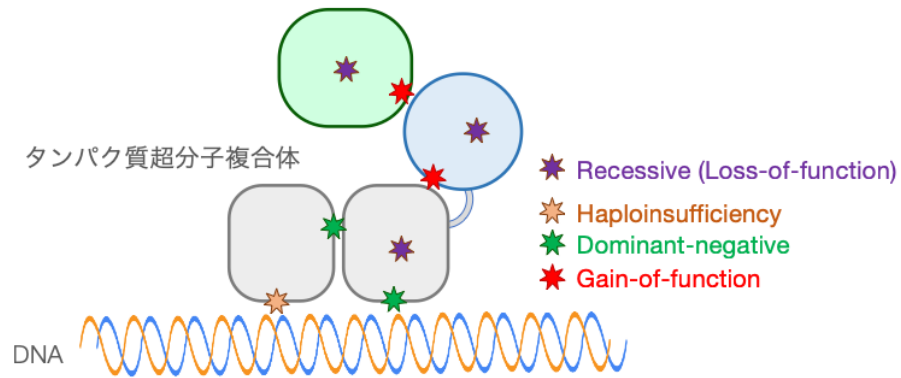


# AD変異は空間的にクラスターになっていることが多い

- AD変異は構造への影響というよりも分子機能を障害することを示唆する
- 同一遺伝子の変異が蓄積している場合はクラスター解析が有効



# タンパク質立体構造から導かれた ゲノムバリエーションと疾患メカニズムとの関係

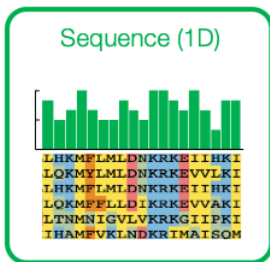


複合体の立体構造情報から疾患表現型を予測できる可能性

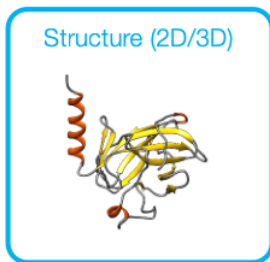


Machine learning approach

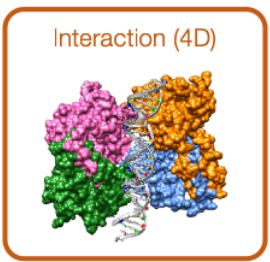
タンパク質複合体の情報は疾患メカニズムの予測に使えるか？



Feature	Value Type
RefAA Frequency	Cont.
MutAA Frequency	Cont.
AA Volume change	Cont.
Grantham	Cont.
PhyloP Mammalian	Cont.
PhyloP Vertebrate	Cont.
PSSM	Cont.
TM-helix	Binary



Feature	Value Type
Solvent accessibility	Cont.
DSSP Coil	Binary
DSSP Strand	Binary
DSSP Helix	Binary
DISOPRED score	Cont.
IDR	Binary



Feature	Value Type
ΔSolvent accessibility	Cont.
HomoSubunit Interface	Binary
HeteroSubunit Interface	Binary
DNA/RNA Interface	Binary
Ligand Interface	Binary
InterDomain Interface	Binary
Distance Protein	Cont.
Distance HomoSubunit	Cont.
Distance HeteroSubunit	Cont.
Distance DNA/RNA	Cont.
Distance Ligand	Cont.

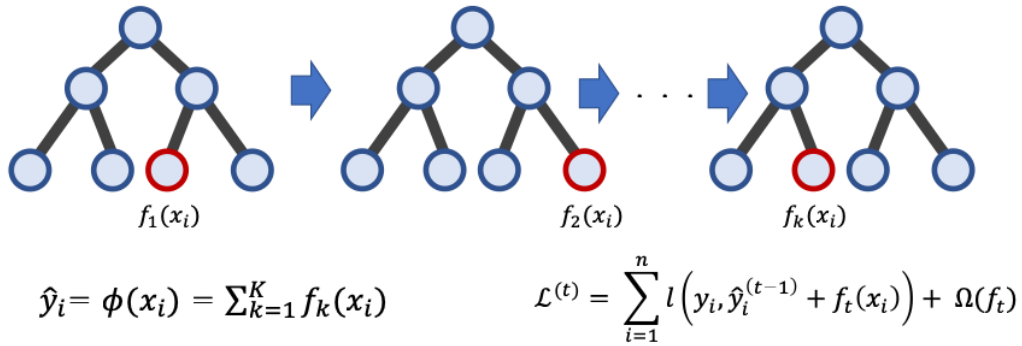


Feature	Value Type
VEST3	Cont.
SIFT	Cont.
FATHMM	Cont.
PROVEAN	Cont.
CADD Phred	Cont.
MutationAssessor	Cont.

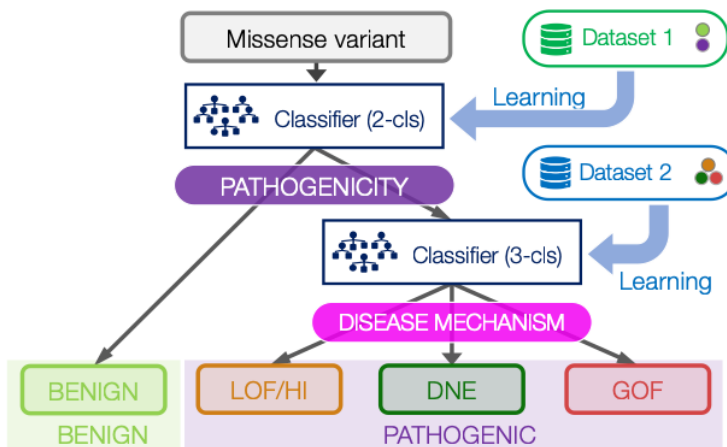


# Gradient Boosting Decision Tree

1. 決定木 (回帰木) を学習させる
2. 1の予測値と目的値から計算される損失関数( $\mathcal{L}$ )が小さくなるように新たな回帰木を学習させ、モデルに追加
3. 2を指定した回帰木の本数分( $K$ )繰り返す
4. モデルの予測値 ( $\hat{y}_i$ )は、データが各回帰木に属する葉の重みの合計となる



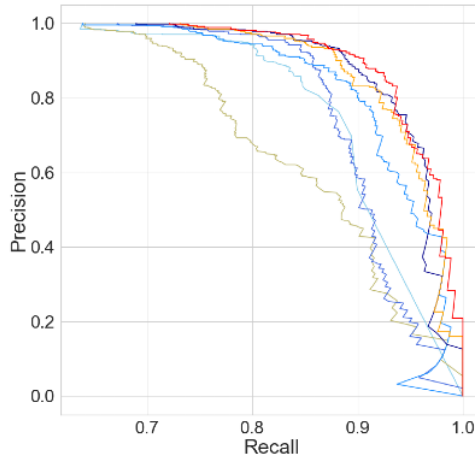
## 疾患メカニズム分類の機械学習モデル



2-Class	3-Class	Training/Test set	
		Genes	Variants
Benign		482	4,786
Pathogenic	LoF/HI	205	1,405
	DNE	173	1,379
	GoF	140	1,200

# 病原性ー良性の2クラス分類の精度評価

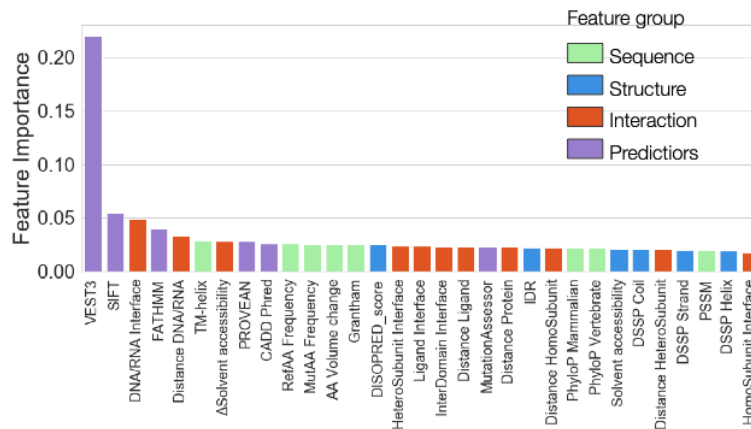
- バリエントが病原性かそうでないかを分類する性能を既存の予測手法と比較
- 既存法でも一部十分な精度が出ているが、相互作用情報を加えた場合が最も高い



Predictor	AP score
SIFT	0.876
FATHMM	0.862
CADD	0.906
VEST3	0.932
REVEL	0.951
	0.951
	<b>0.960</b>

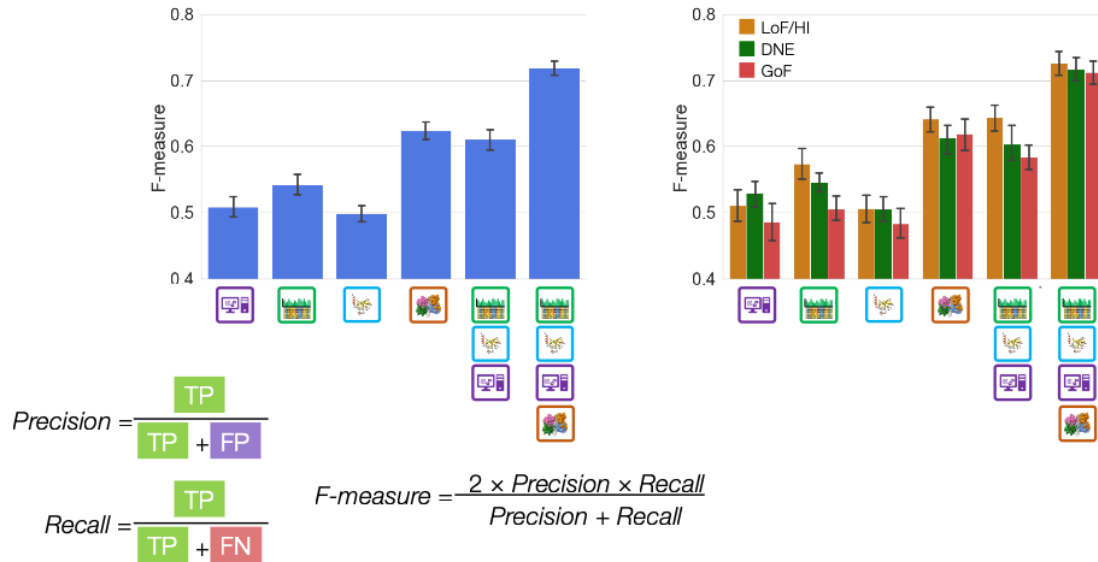
# 2クラス分類における特徴量の重要度

- 病原性ー良性の分類では既存の予測法のスコアが最も分類性能に寄与
- ある意味でこの分類手法については概ね完成しているといえる



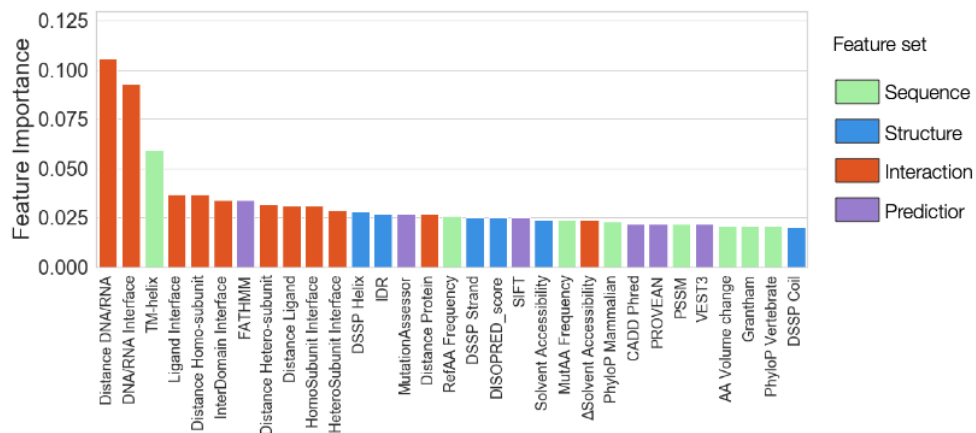
## バリエーションの疾患メカニズムの分類性能の評価

- 4種類の特徴量セットの組み合わせでバリエーションの疾患メカニズムの分類性能を評価
- 分子間相互作用の特徴量を用いた場合に分類性能が飛躍的に向上した



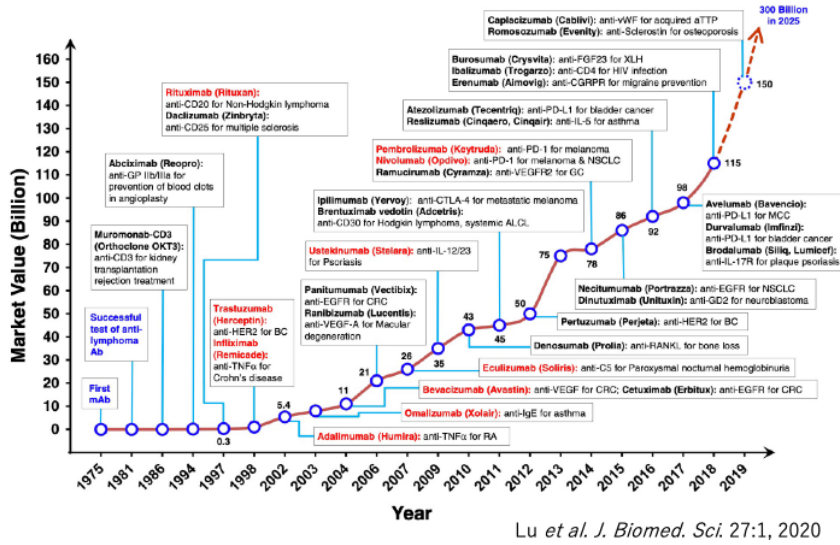
## どの特徴量が分類性能に影響を与えているか？

- 疾患メカニズム(LoF, DNE, GoF)の分類においては、特にDNA/RNA相互作用に関する特徴量が重要度が高い
- 総じて相互作用に関する特徴量が重要度が高い傾向にあった
- 既存の予測手法は分類性能にほとんど貢献していない (あってもごくわずか)



# 抗体医薬品の市場規模は年々拡大している

- 最初期はがんの治療薬として開発。現在はがん以外の治療薬としても重要性が増している。
- 570を超える抗体医薬品が治験中あるいは既承認。市場規模は2025年に3000億ドルを超える見込み。



免疫チェックポイント阻害薬  
オプジーボ (Nivolumab)

## 2021年度 医療用医薬品 国内売上高トップ20

【単位：億円、%】

順位	製品名	薬効・領域	社名	21年度売上高		
				前年比	前年比	
Pembrolizumab Nivolumab	1	キイトルーダ*	がん	MSD	1195	1.1
	2	オプジーボ	がん	小野薬品工業	1124	13.8
	3	タグリソン*	がん	アストラゼネカ	1037	9.1
	4	タケキャブ	消化性潰瘍	武田薬品工業	946	13.6
Bevacizumab Casirivimab/Imdevimab	5	リクシアナ	抗凝固薬	第一三共	925	19.5
	6	ネキシウム*	抗胃酸薬	アストラゼネカ	913	▲ 2.0
	7	イグザレルト*	抗凝固薬	バイエル薬品	815	2.9
	8	アバステン	がん	中外製薬	809	▲ 0.7
Atezolizumab Ramucirumab Ustekinumab	9	ロナプリーブ	感染症	中外製薬	774	—
	10	アジルバ	高血圧症	武田薬品工業	763	▲ 7.2
	11	サムスカ	利尿薬	大塚製薬	730	11.8
	12	アイリリア	加齢黄斑変性など	参天製薬	725	12.5
Adalimumab Golimumab	13	デセントリク	がん	中外製薬	622	65.9
	14	サイラムザ*	がん	日本イーライリリー	526	0.9
	15	ステララ	乾癬・潰瘍性大腸炎など	田辺三菱製薬	515	59.9
	16	オフエブ*	抗緑膿菌薬	NBI	510	29.8
	17	ヒュミラ	リウマチなど	エーザイ	506	▲ 2.5
	18	サインバルタ*	抗うつ薬	日本イーライリリー	501	▲ 21.4
	19	イクスタンジ	がん	アステラス製薬	472	17.3
	20	シンボニー	リウマチなど	田辺三菱製薬	433	2.4



# モノクローナル抗体命名ルール (旧\*)

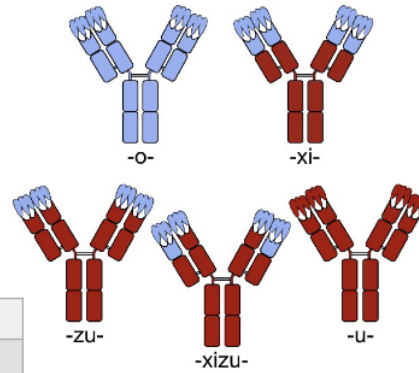
<プレフィックス>-<ターゲットサブシステム>-<ソースサブシステム>-mab

ターゲットサブシステム

-t(u)-	腫瘍	<b>tumor</b>
-l(i)-	免疫調節	<b>immunomodulating</b>
-ci-	心血管	<b>cardiovascular</b>
-ne-	神経系	<b>neural</b>
-vi-	ウイルス	<b>viral</b>

ソースサブシステム

-mo-	マウス抗体	<b>mouse Ig</b>
-xi-	キメラ型抗体	<b>chimeric Ig</b>
-zu-	ヒト化抗体	<b>humanized Ig</b>
-xizu-	キメラ型/ヒト化抗体	<b>chimeric/humanized Ig</b>
-u-	ヒト抗体	<b>human Ig</b>



Nivo-l-u-mab

Ce-tu-xi-mab

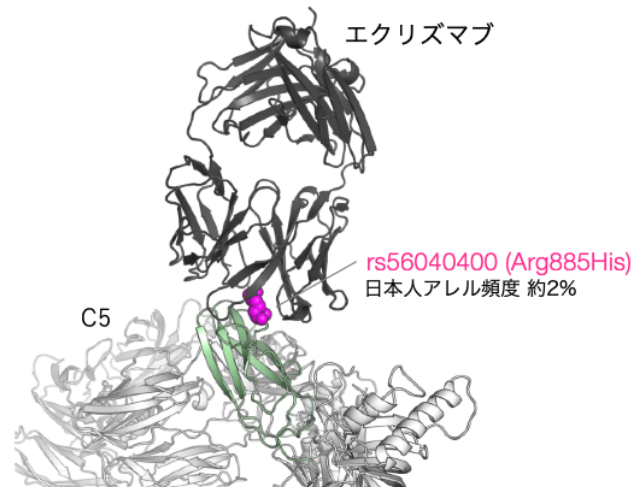
Beva-ci-zu-mab

\* 2021年からは命名ルールが変更されている

## 抗体医薬品ターゲットのゲノムバリエーションによる結合親和性の違い？

- 医薬品の効果は人によって異なっており予測が困難なことが多い
- ゲノムビッグデータから各バリエーションの薬剤応答性を予測できるか？

- エクリズマブ（一般名ソリリス）は、発作性夜間へモグロビン尿症の治療薬（抗体医薬品）。
- 日本の治験で、エクリズマブへの反応が悪い患者が多いことが報告され、調べたところ効果の低い患者に共通のバリエーション（R885H）が見つかった（Nishimura et al. 2014）。
- エクリズマブとC5の複合体構造において、R885Hは、エクリズマブとの相互作用面にあり、結合を阻害していると考えられる（Asbjørn et al. 2016）。
- 現在では、この部分をエピトープとしない抗体医薬品（Crovalimab）も開発されている。
- その他のレアバリエーションについては影響は不明。

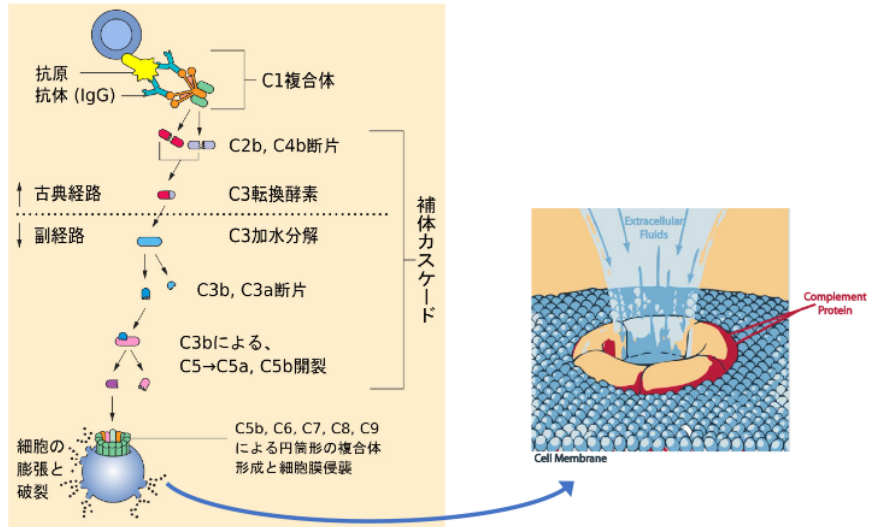


薬剤応答性を正しく予測できれば治療薬の選択・精密医療が可能

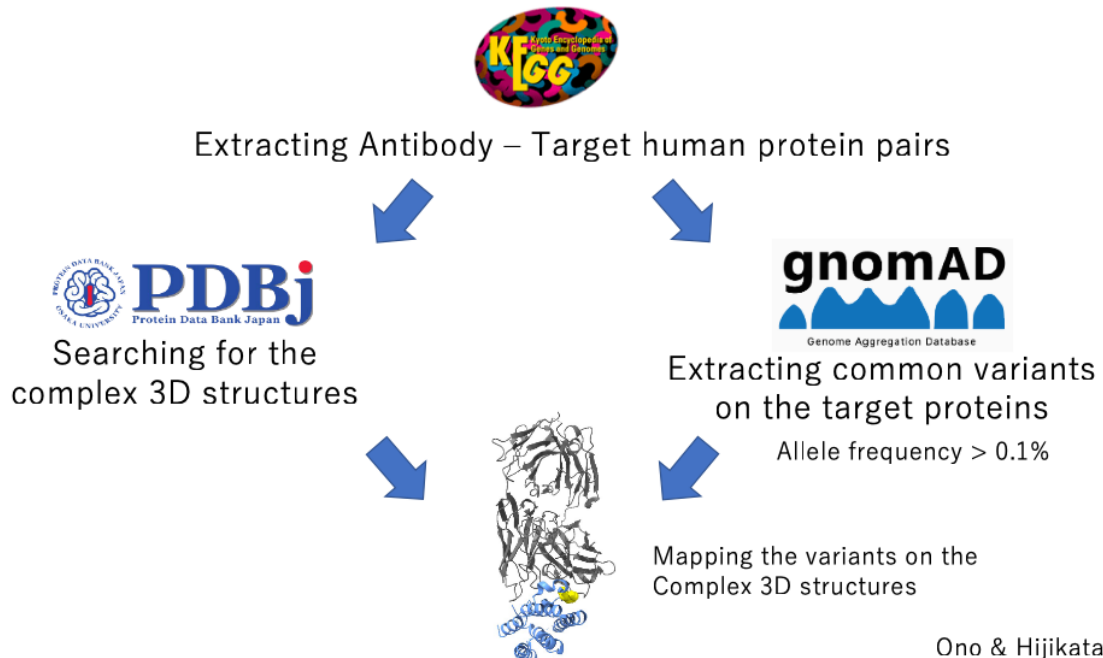


# 補体カスケード

- 免疫系のタンパク質群で、抗体やマクロファージを補助する役割を持つ
- PNHでは、補体系が赤血球を攻撃することで溶血が起こる



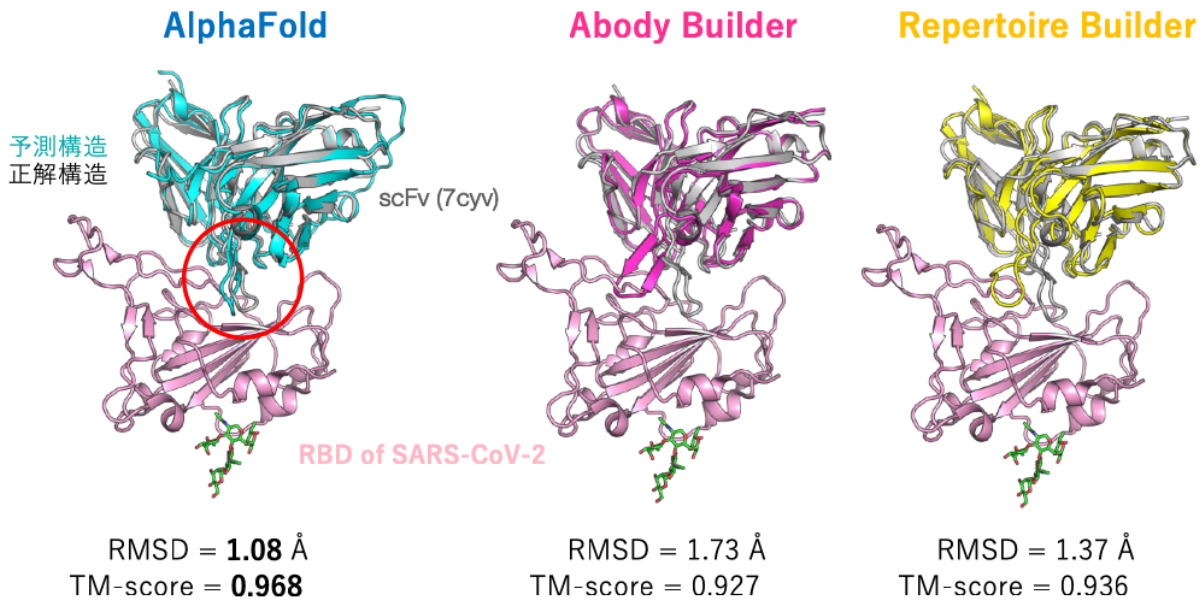
抗体医薬品ターゲットの構造データ及びゲノムバリエーションデータの取得





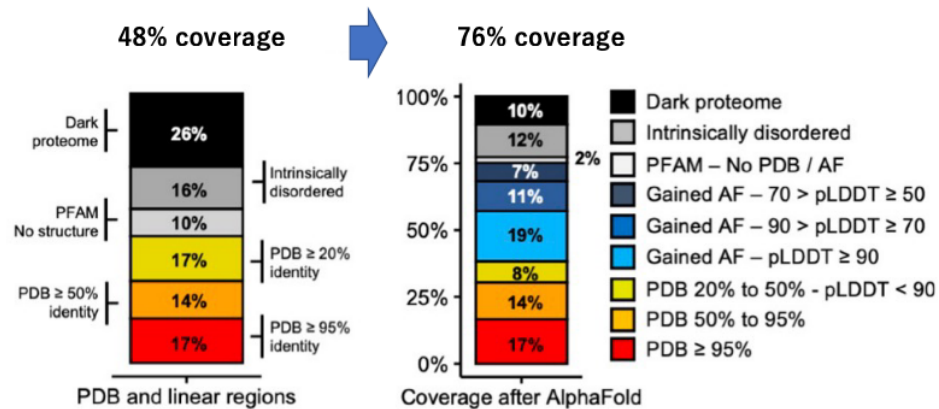
# AlphaFoldは構造未知の抗体構造予測に使えるか？

- AF2を用いた抗体の立体構造のブラインドテストの結果の一例
- 既存法よりも精度が高い、特にCDR領域(抗原認識部位)の予測性能が高い



## AF2によってヒトプロテオームはどのくらい拡張されたか

- AF2によって、ホモロジーモデル可能なものも含めた構造48%だった構造カバレッジは76%まで増大
- なんの構造の手がかりもなかった「Dark proteome」も26%→10%まで減少
- ただし単量体かつ低分子結合情報はない



# AF2モデルから複合体構造を精度よく予測できつつある

ARTICLE

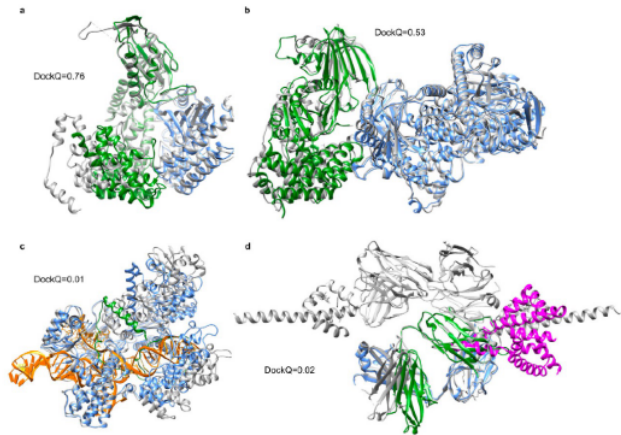
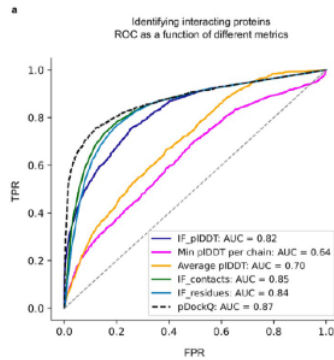
Check for updates

<https://doi.org/10.1038/s41467-022-28865-w>

OPEN

## Improved prediction of protein-protein interactions using AlphaFold2

Patrick Bryant <sup>1,2,3✉</sup>, Gabriele Pozzati <sup>1,2,3</sup> & Arne Elofsson <sup>1,2✉</sup>



Bryant et al. *Nat Commun.* (2022)

# モノクローナル抗体命名ルール

旧ルール(~2020)

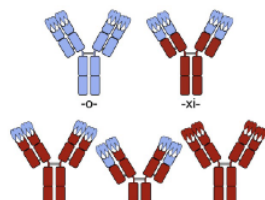
-momab	マウス抗体	<b>M</b> ouse mAb
-ximab	キメラ型抗体	<b>Ch</b> imeric mAb
-zumab	ヒト化抗体	<b>H</b> umanized mAb
-umab	ヒト抗体	<b>H</b> uman mAb

-t(u)-	腫瘍	Tumor
-l(i)-	免疫調節	Immunomodulating
-ci-	心血管	Cardiovascular
-ne-	神経系	Neural
-vi-	ウイルス	Viral

新ルール(2021~)

-tug	未修飾抗体	<b>U</b> nmodified immunoglobulins
-bart	人工抗体	Antibody <b>art</b> ificial
-mig	多重特異性抗体	<b>M</b> ulti-immunoglobulin
-ment	フラグメント抗体	Fragment

-ler-	アレルゲン	<b>All</b> ergen
-pru-	免疫抑制	Immunosup <b>pr</b> essive
-sto-	免疫賦活	Immunost <b>imulato</b> ry



Nivo-l-umab

Ce-tu-ximab

Leca-ne-mab

## まとめと展望

- ゲノムバリアントの表現型(疾患・薬の応答性)への影響を評価する上で、分子構造情報(特に分子間相互作用)が有効
- ゲノム情報に比べて、構造情報はまだまだ不足しているのが現状
- 疾患変異の分子メカニズムの情報も限定的である(データベースの整備が必要)
- 構造生物学の進展とAIを活用した「3Dインタラクティブ」が鍵