

「PDBjタンパク質をゲノムに マップしたpdbBAMの作成」

城田松之

東北大学大学院医学系研究科

創生応用医学研究センター

新医学領域創生分野

平成27年3月28日(土)

平成26年度「統合データ解析トライアル」研究成果報告会

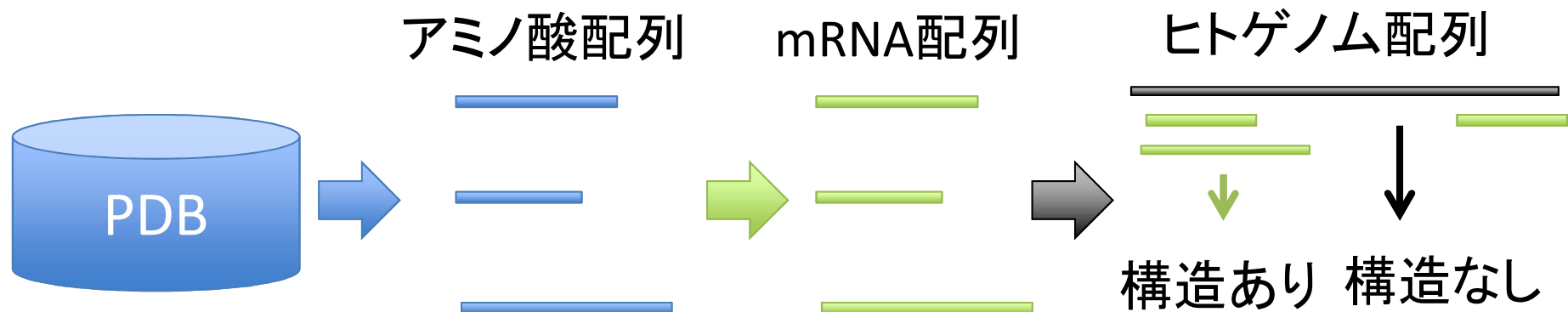


本日の発表内容

- pdbBAMとは
- ツールの作成
- IGVを用いた個人ゲノム情報とPDB情報の表示
- まとめ

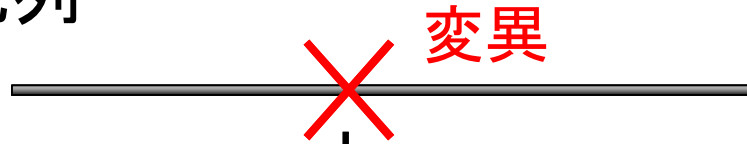
pdbBAMとは何か

- PDBに含まれるタンパク質のアミノ酸配列をmRNA配列を媒介してゲノムにマップしてBAM形式としたもの
- PDB全体をゲノム全体に貼付ける
- ゲノムのどこの遺伝子が構造解析されているかを一目でわかるようにする



これまでの構造情報の利用

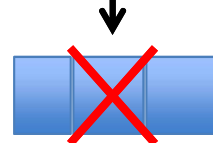
ゲノム配列



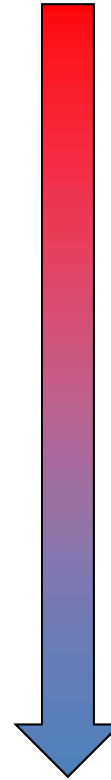
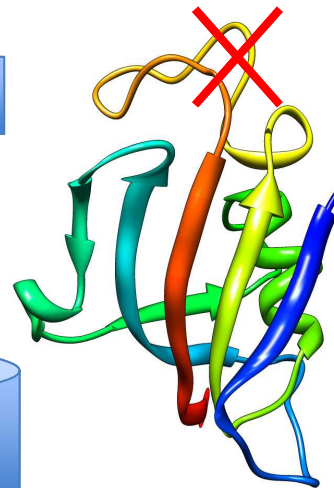
mRNA配列



タンパク質配列

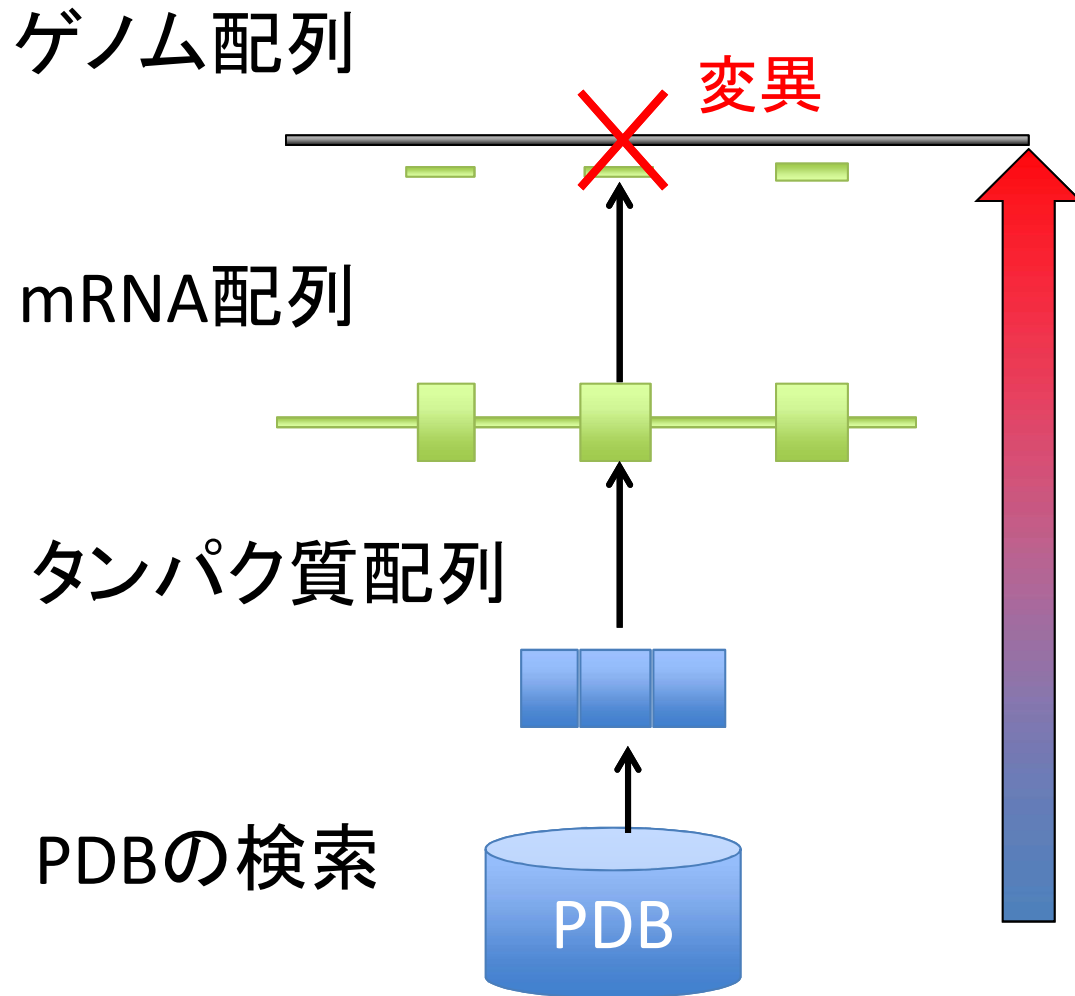


PDBの検索



- StSNP
- coliSNP
- MuPIT
- など

PDB配列からゲノム配列への逆マッピング



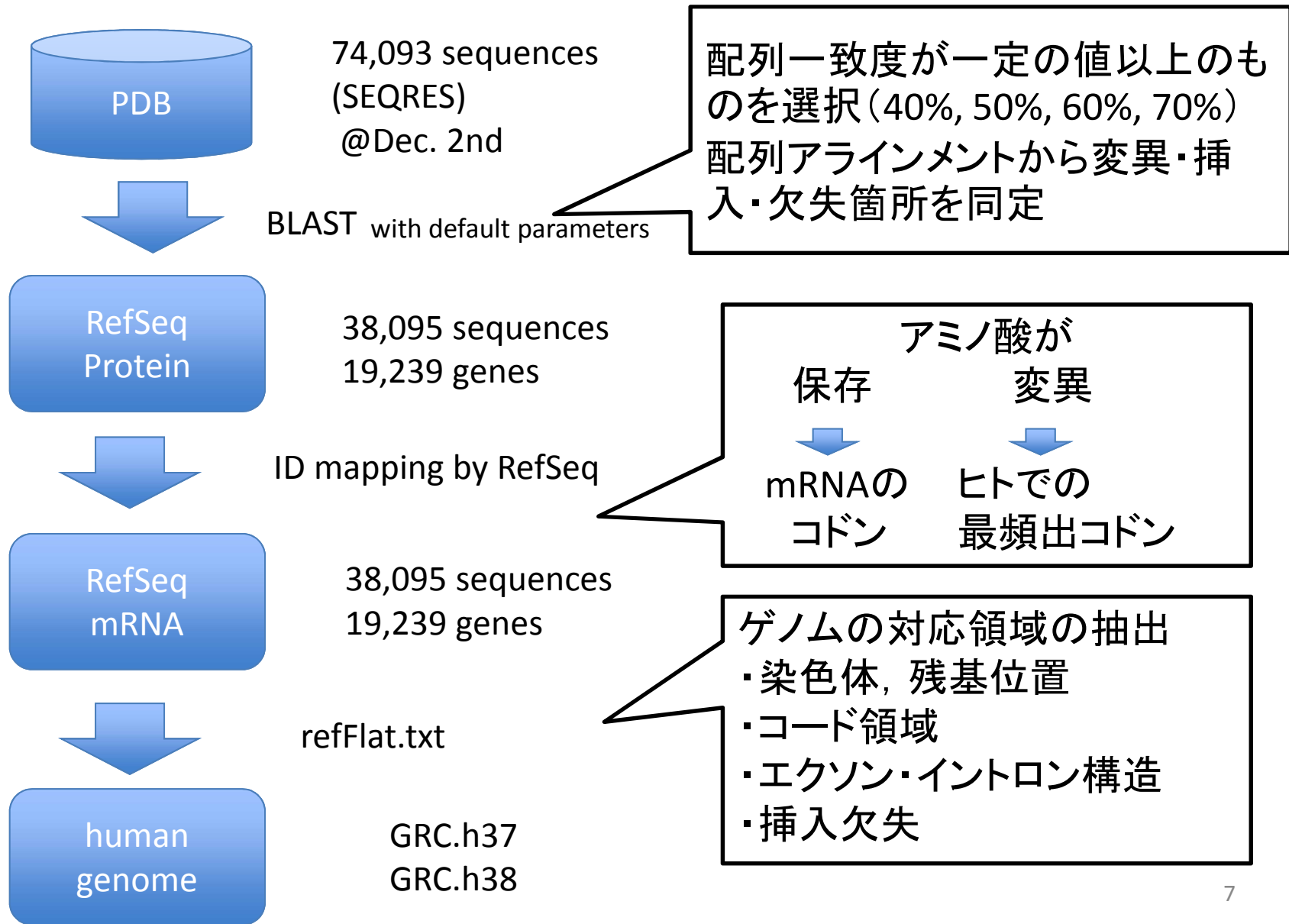
- 利点
 - 変異ごとの検索なし
でよい
 - ゲノムを見ればPDB
全体が分かる
- 欠点
 - 逆翻訳は曖昧さが
残る
 - PDB検索の閾値

この向きのマッピングを行うツールはまだ存在しない

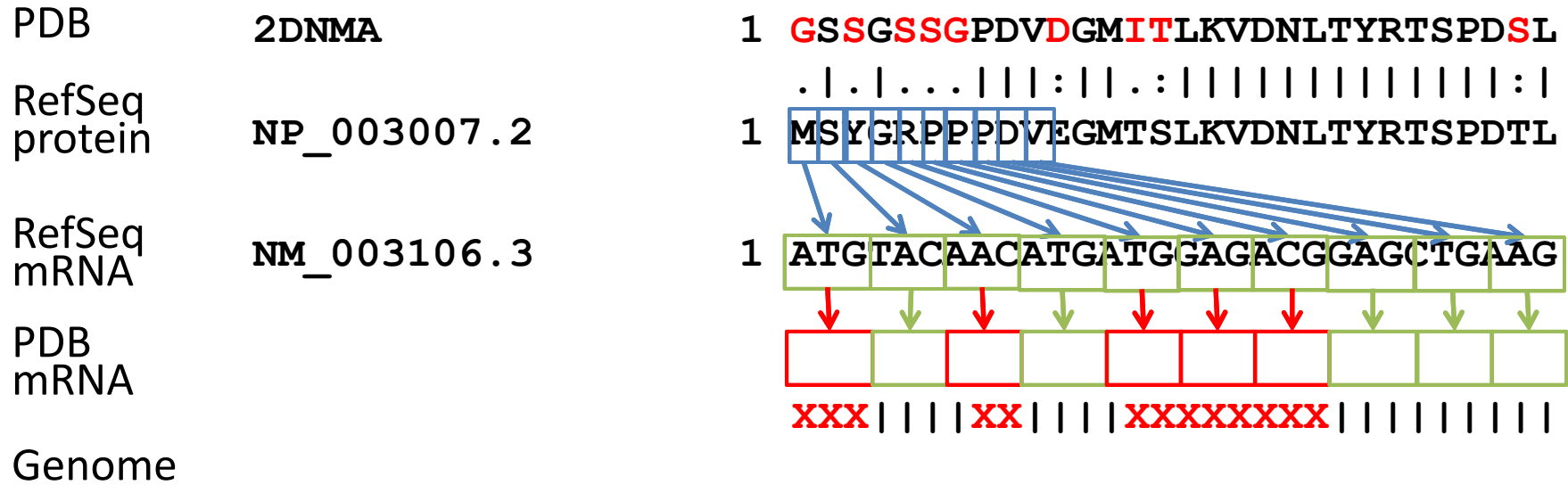
利用したデータベース

- 日本蛋白質構造データバンク(PDBj)
 - タンパク質の立体構造およびアミノ酸配列情報
- NCBI RefSeq
 - タンパク質アミノ酸配列と対応するmRNA配列
- Genome Reference Consortium (GRC)
 - ヒトゲノム配列
 - GRC.h38(最新版)とGRC.h37(1つ前の版)

pdbBAM作成の流れ



PDB配列からゲノムへのアラインメント

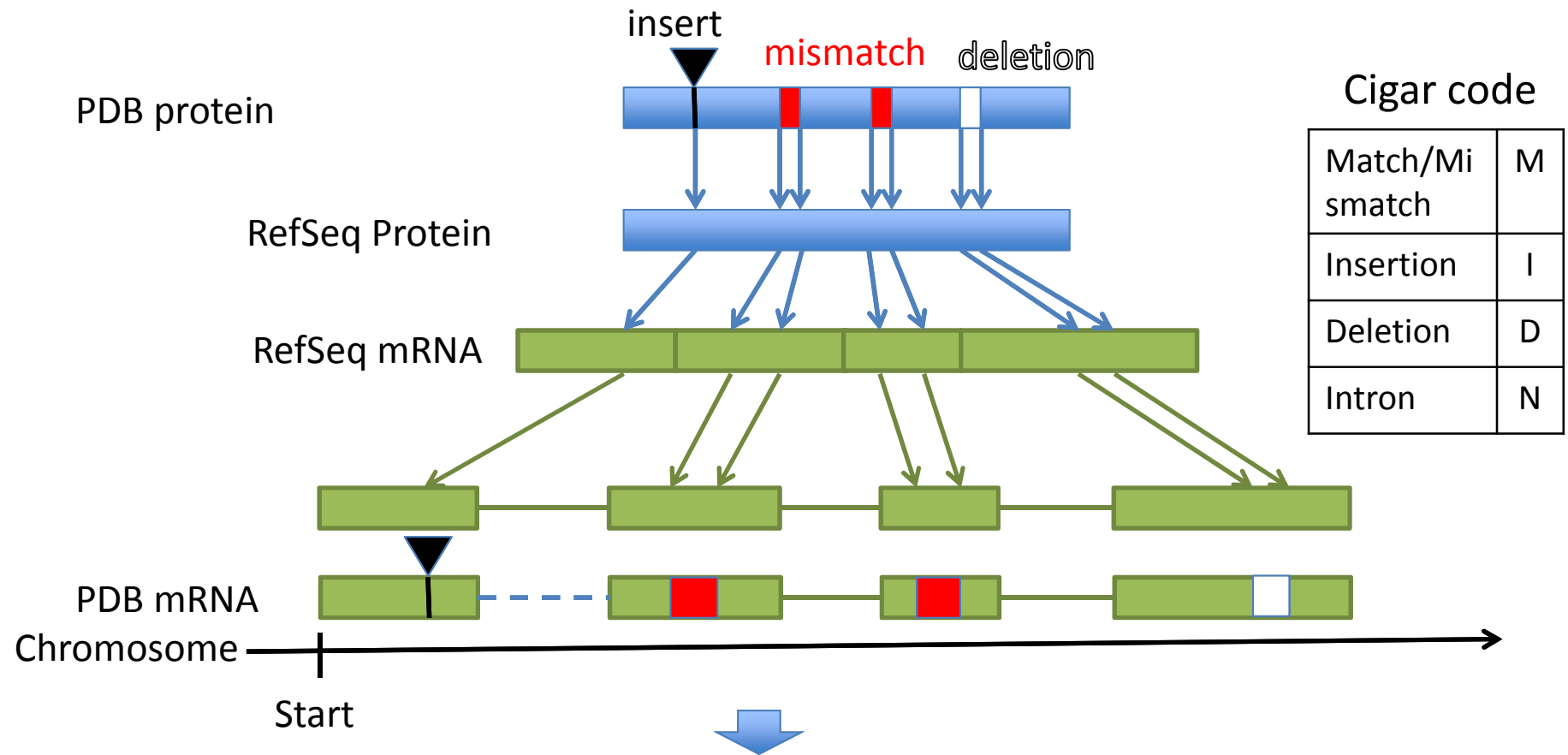


mRNAと一致 → mRNAのコドンを利用

mRNAと不一致 → アミノ酸に対する最頻出コドンで置換

アミノ酸	コドン	アミノ酸	コドン	アミノ酸	コドン	アミノ酸	コドン	アミノ酸	コドン
A	GCC	C	TGC	D	GAC	E	GAG	F	TTC
G	GGC	H	CAC	I	ATC	K	AAG	L	CTG
M	ATG	N	AAC	P	CCC	Q	CAG	R	AGA
S	AGC	T	ACC	V	CTG	W	TGG	Y	TAC

アラインメントからSAMフォーマットへ



SAMフォーマット

```
pdb|3WHD|A_1 0 chr12 8670819 255 52M734N152M1065N116M782N145M * 0 0 CATGCA...
```

配列名

参照配列 位置

Cigar code

塩基配列₉

SAM/BAM変換

SAM -> BAM変換

```
samtools view -Sb pdbbam.sam | samtools sort > pdbbam.bam
```

インデックス作成

```
samtools index pdbbam.bam
```

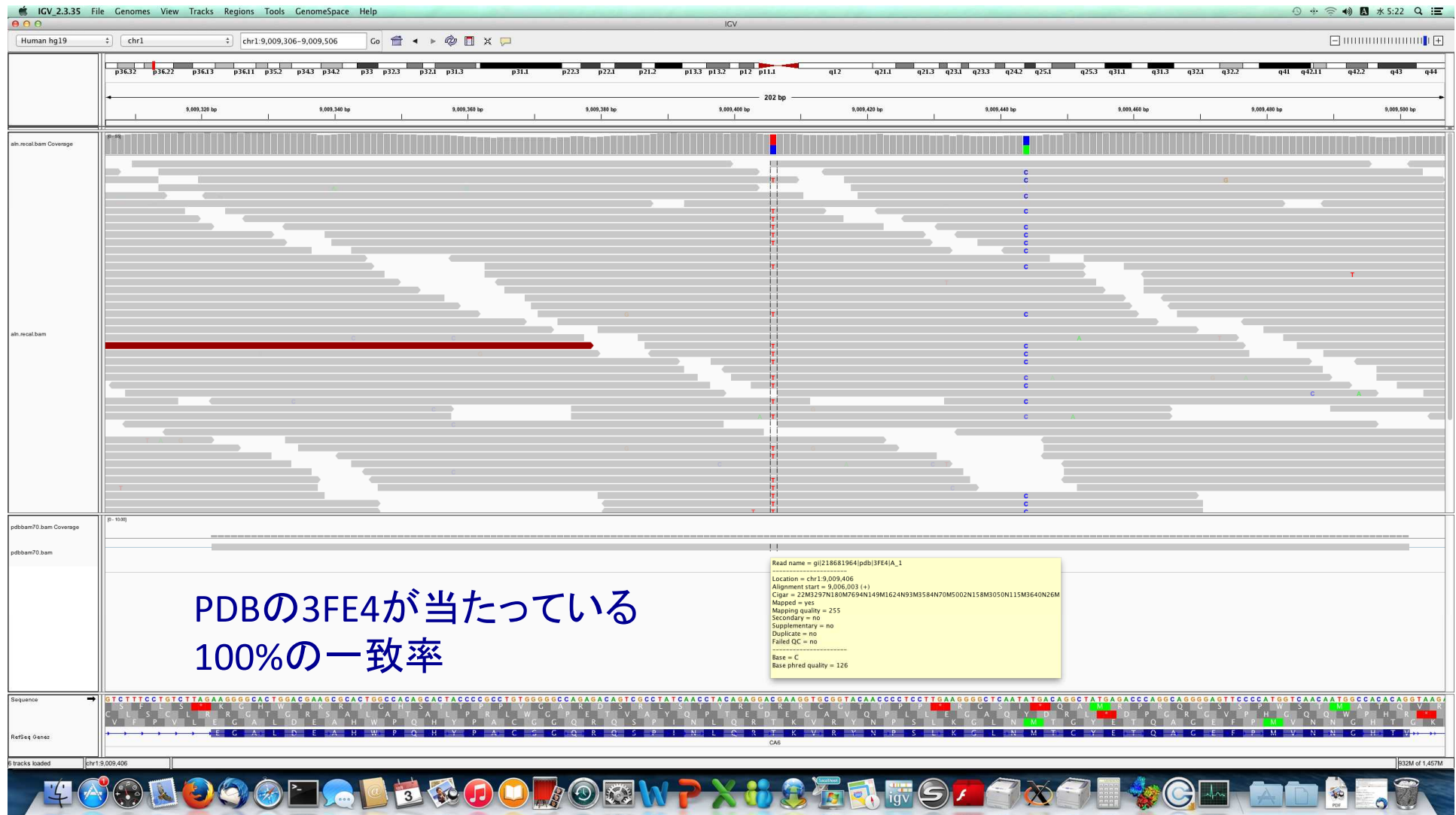
エクソン情報にあるPDBの抽出

```
bedtools coverage -hist -abam pdbbam.bam -b refFlat.bed >  
pdbbam.coverage.txt
```

samtools 0.1.18 Li H et al. Bioinformatics 25(16):2078-9, 2009

BEDtools 2.150.0 Quinlan AR, Hall IM. Bioinformatics 26(6):841-2, 2010

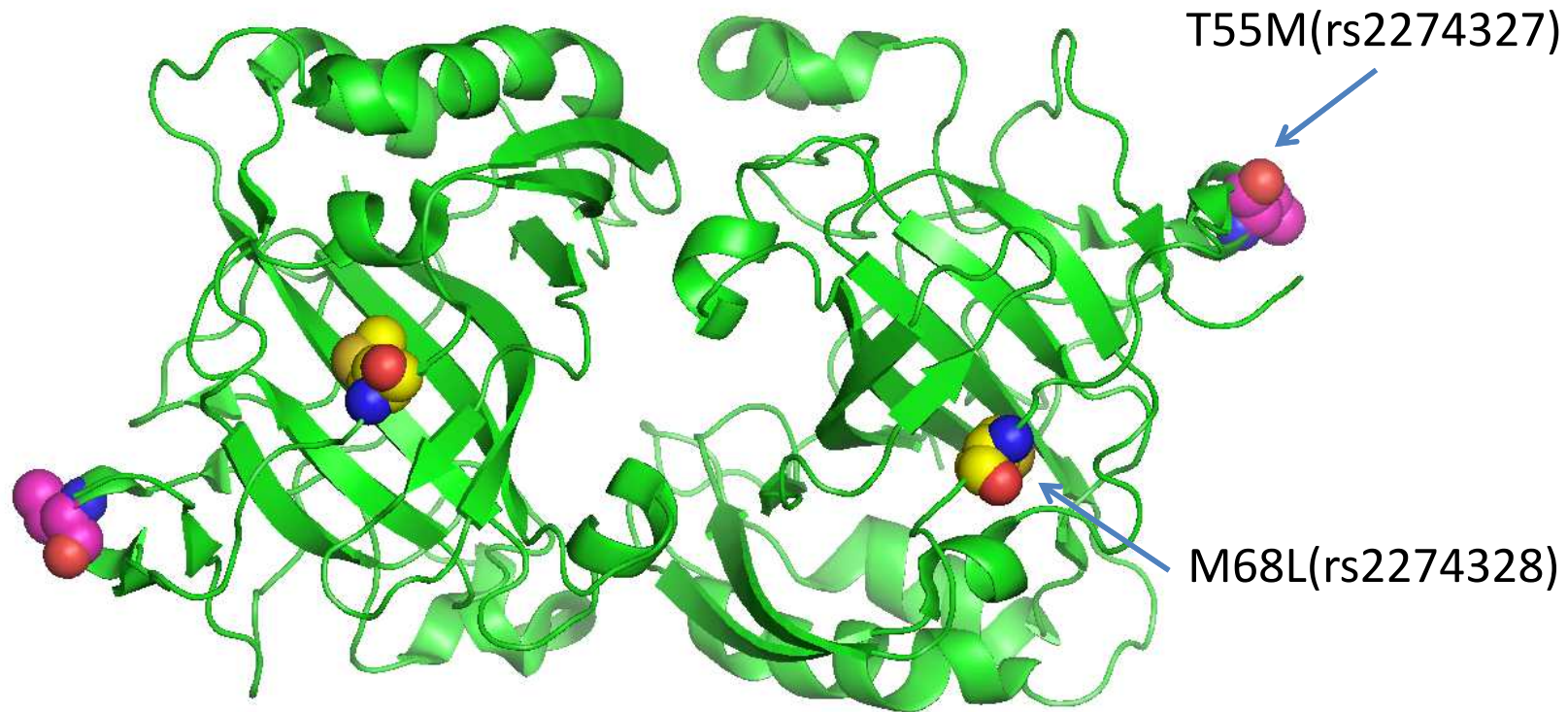
実際の例 (Carbonic Anhydrase遺伝子)



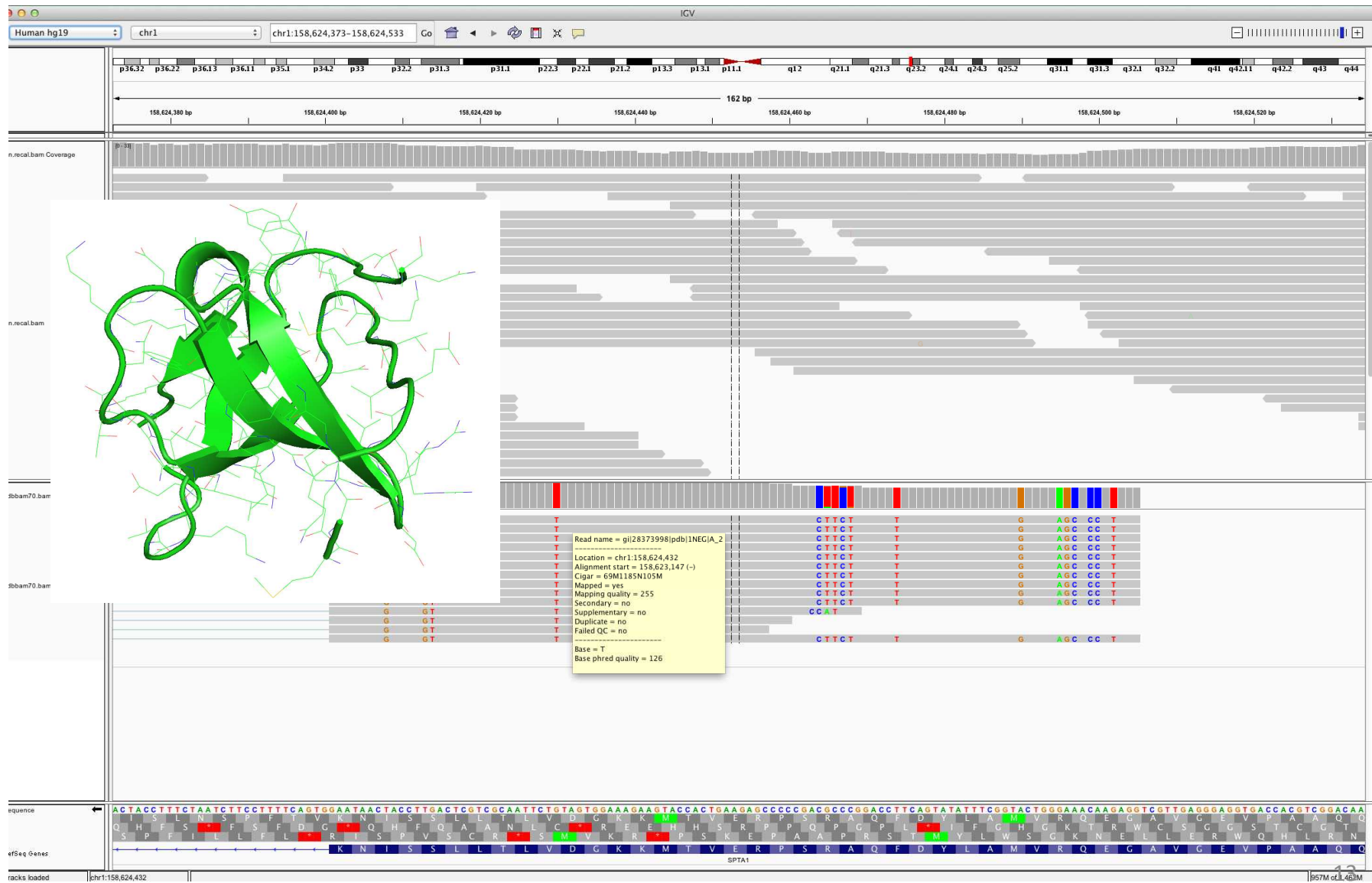
PDBの3FE4が当たっている
100%の一致率

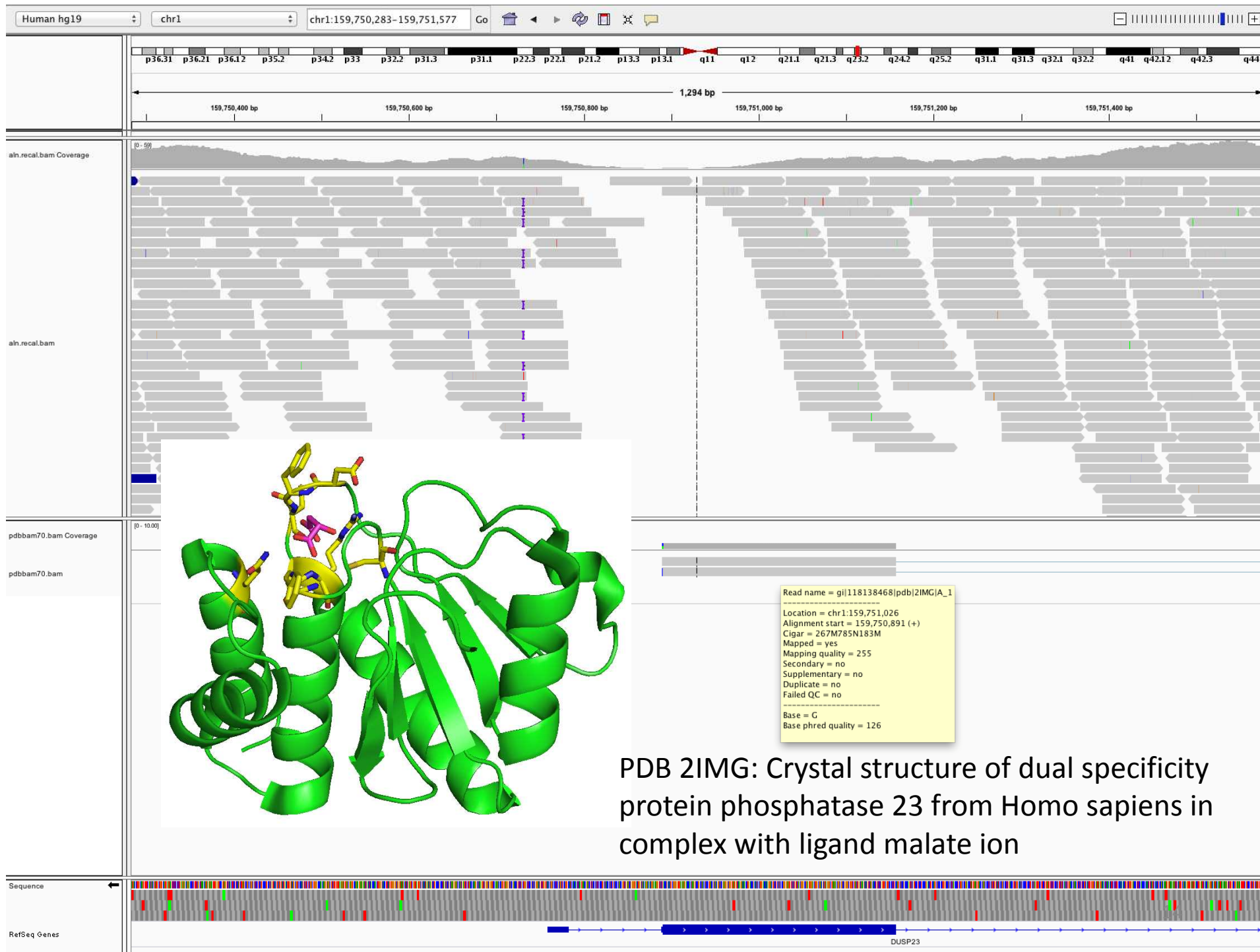
上段: NGSによる個人ゲノム解析結果
中断: pdbBAM70
下段: 遺伝子領域と翻訳

Carbonic Anhydrase (CA6) PDB ID 3FE4



SH3 and cysteine domain-containing protein isoform 1 (SPTA1)





P53領域



上段: NGSによる個人ゲノム解析結果
中断: pdbBAM70
下段: 遺伝子領域と翻訳

pdbBAMまとめ

- PDBのタンパク質を配列相同性を用いてヒトゲノムにマッピングしてpdbBAMを作成
- 配列一致度について40%~70%までで作成
- ヒトゲノム上のコード領域におけるカバー率と平均深度は
 - カバー率: 19%(配列一致度70%) ~ 24%(同40%)
 - 平均深度: 0.8(70%) ~ 2.4(40%)
- 個人ゲノム解析結果とあわせてゲノム上のどこに、どの程度の数と類似度の立体構造情報があるかを表示できる