

博士論文

三次元分布関数を用いた  
タンパク質周辺の水分子のサンプリング手法の開発

北里大学 薬学部

小林慎平

# 目次

第 1 章. 序論.....	3
1-1. 背景 .....	3
1-2. 結晶水と予測水 .....	3
1-3. 水の三次元分布関数 .....	4
1-3-1. 概要 .....	4
1-3-2. 3D-RISM 理論 .....	5
1-3-3. 三次元分布関数を基にした水分子の予測 .....	6
1-3-4. Placevent .....	7
1-3-5. 本研究の目的 .....	7
第 2 章. 方法.....	9
2-1. 本研究における重み付きモンテカルロ法 .....	9
2-2. 水分子のサンプリングに対する重み付きモンテカルロ法の適用 .....	10
2-3. 調整重みの導入 .....	11
2-4. プログラム構築 .....	13
2-4-1. 予測水同士の距離 .....	13
2-4-2. DroPred のアルゴリズム .....	15
2-4-3. 並列化による DroPred 計算の効率化 .....	16
2-5. 性能評価 .....	17
2-5-1. 性能評価プロトコル .....	17
2-5-2. テストセットの準備 .....	18
2-5-3. 水の三次元分布関数の計算 (3D-RISM 計算) .....	18
2-5-4. 水分子の位置予測 .....	20
2-5-5. DroPred の性能評価 .....	21
第 3 章. 結果・考察.....	22
3-1. 水分子の再現の定義 .....	22
3-2. DroPred の出力サンプル数 .....	23
3-3. 予測水の配置数 .....	24
3-4. 評価 1: 結晶水の位置の再現度による評価 .....	25
3-4-1. 評価の概要 .....	25

3-4-2. Coverage の定義 .....	27
3-4-3. ターゲット毎の結果.....	28
3-4-4. テストセット全体の傾向 .....	32
3-5. 評価 2：分布関数の再現度による評価.....	33
3-5-1. 評価の概要 .....	33
3-5-2. 結晶水が存在する場所の確率密度 .....	33
3-5-3. 予測水の出現率 .....	34
3-5-4. 分布関数の再現度についての結果 .....	35
3-6. 配置数の制限の影響 .....	37
3-7. Placevent と DroPred の比較考察.....	38
第 4 章. 今後の展望.....	41
第 5 章. 結論.....	43
謝辞.....	44
参考文献 .....	45
論文目録 .....	47
補足資料 .....	49
S-1. 検証用テストセット .....	49
S-2. 各手法における水分子の配置数 .....	54
S-3. ターゲット毎の検証結果.....	59
S-4. Coverage の中央値 .....	79

# 第1章. 序論

## 1-1. 背景

タンパク質は、生体内において多数のアミノ酸が鎖状に連なった状態で生成され、折りたたまれて二次構造・三次構造を形成することでその機能を発現する。また、多くのタンパク質が他のタンパク質やリガンドなどと結合し、相互作用することも知られている。タンパク質の構造形成・機能の発現や、他のタンパク質、リガンドとの相互作用において、タンパク質周辺に存在する水分子の存在はきわめて重要な要素である。このため、タンパク質周辺の水分子の位置を知ることが、タンパク質の機能や相互作用の理解、ドラッグデザインなどの創薬研究などにつながると言える。

現在、タンパク質周辺の水分子の位置は、X線結晶構造解析、クライオ電子顕微鏡などの実験的な手法によって解析されている。実際に、実験によって解析されたタンパク質の構造情報を収集しているデータベースである Protein Data Bank (PDB)<sup>1</sup>に登録されているタンパク質の立体構造データの多くには、タンパク質周辺に存在する水分子の位置情報が含まれている。

また、タンパク質周辺の水分子の位置は、コンピュータ計算を用いた手法によって予測することも可能である<sup>2-9</sup>。コンピュータ計算を用いた代表的な予測手法としては、MD (分子動力学) シミュレーションや、水分子のドッキング、統計処理などを行う知識ベースのアプローチなどが挙げられる。コンピュータを用いた予測手法は、近年のコンピュータ関連の技術の発達もあり、実験的な手法と比較して、小さな時間的・費用的なコストで実行することが出来る可能性がある。またタンパク質の立体構造予測の分野では近年、「AlphaFold」を始めとする、高精度の立体構造予測手法も使用可能になっている<sup>10,11</sup>。コンピュータ計算による水分子の位置予測は、予測手法によって予測された立体構造 (モデル構造) など、仮想的なタンパク質にも適用可能であるという利点もある。

このような背景から、タンパク質周辺の水分子の位置を高精度で予測する手法を開発することができれば、タンパク質の構造や機能の解析や、ドラッグデザインなどの研究に貢献することが出来ると考えられる。本研究室においても、タンパク質周辺の水分子の位置予測手法の開発に取り組んでいる。

## 1-2. 結晶水と予測水

本論文においては、実験手法を用いてその位置が解析されたタンパク質周辺の水分子を「結晶水」、コンピュータ計算を用いた予測手法によってその位置が

予測された水分子を「予測水」とそれぞれ呼ぶこととする。また、予測手法によってあるタンパク質周辺における水分子の位置を予測した1つの水和構造を「サンプル」と呼ぶこととする。

### 1-3. 水の三次元分布関数

#### 1-3-1. 概要

「水の三次元分布関数」の計算は、タンパク質周辺における水分子の位置を予測するためのアプローチの1つとして知られている。この手法は、タンパク質などの溶質の影響を全く受けない場所（バルク）における水分子の存在しやすさを1と定義した場合の、計算領域内の各位置における水分子の存在しやすさを計算するものである。三次元分布関数によって与えられる各位置の水分子の存在しやすさは「水の確率密度」と呼ばれる。例えば、ある位置における水の確率密度が5.0である場合、その位置にはバルクの5.0倍水分子が存在しやすいことになる。タンパク質周辺の水分子の位置を予測する手法には、水分子の位置を直接予測、即ち水分子が存在すると考えられる座標を出力するものも存在するが、三次元分布関数はそれらの手法とは異なり、計算領域全体について水分子の存在しやすさが計算され、連続関数として出力されるという特徴がある。

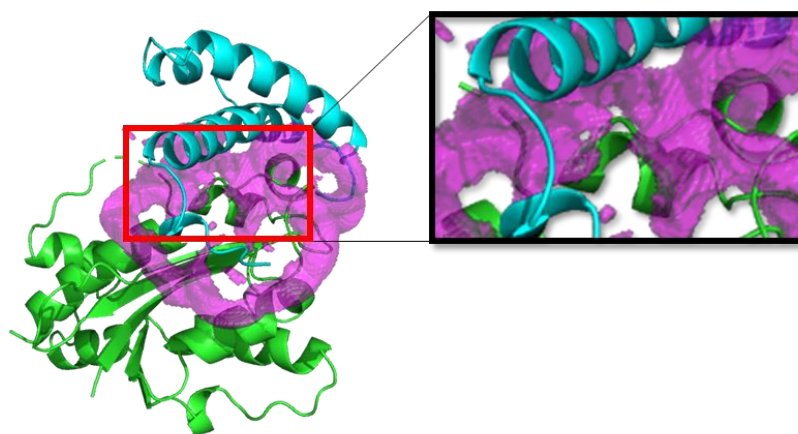


図1. タンパク質-タンパク質界面における水の三次元分布関数の例  
緑色・水色で示したタンパク質-タンパク質複合体 (PDB ID: 2ID0) の界面における水の確率密度が0より大きな領域を紫色の透過面で示している<sup>29</sup>。

図1は、タンパク質-タンパク質複合体の界面における三次元分布関数の例を示したものであり、水分子の存在しやすさが0より大きいと予測された領域を紫の透過面で示している。

水の三次元分布関数は3D-RISM理論<sup>12-15</sup>を用いた計算によって計算されるほか、近年ではAIを用いて三次元分布関数を計算する手法も開発されており<sup>16</sup>、今後三次元分布関数によって水分子の位置を予測するアプローチは更に興味深いものとなる可能性がある。

### 1-3-2. 3D-RISM 理論

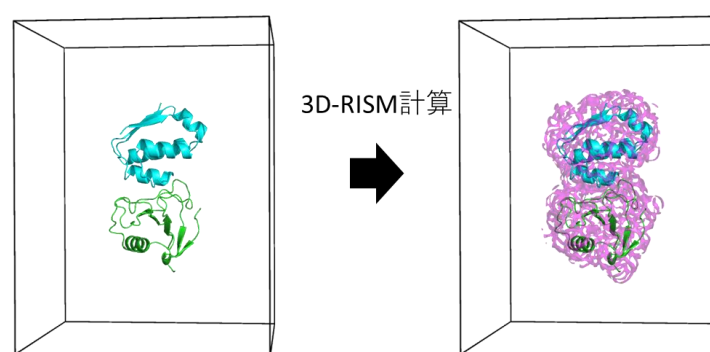


図2. 3D-RISM 理論による三次元分布関数の計算

3D-RISM 理論による計算で水の三次元分布関数を計算することができる。図は、PDB ID: 1AY7 について 3D-RISM 理論で計算した三次元分布関数を示している。この図では、確率密度が 3.0 以上の領域を紫色の透過面で示している。

水の三次元分布関数を計算するための代表的な手法の1つが、3D-RISM(three-dimensional reference interaction site model)理論による計算である。3D-RISM 理論は、統計力学理論の1つであり、3D-RISM 方程式と、その Closure 方程式 (Kovalenko-Hirata Closure 方程式などが代表的な Closure 方程式として挙げられる) を解くことで、タンパク質周辺の水の三次元分布関数を求めるものである。3D-RISM 理論による計算は、水分子の位置を予測するタンパク質 (以降「ターゲット」と呼ぶ) をコンピュータ上の三次元グリッド空間上に配置し、ターゲットを中心とした十分な大きさの直方体のボックスを計算領域

として行う (図 2)。水分子の存在しやすさ (確率密度) は、このボックス内の全てのグリッド点について計算される。尚、三次元グリッド空間におけるグリッド点同士の間隔 (グリッド幅) は、3D-RISM 計算の実行時に自由に設定することが可能である。

### 1-3-3. 三次元分布関数を基にした水分子の予測

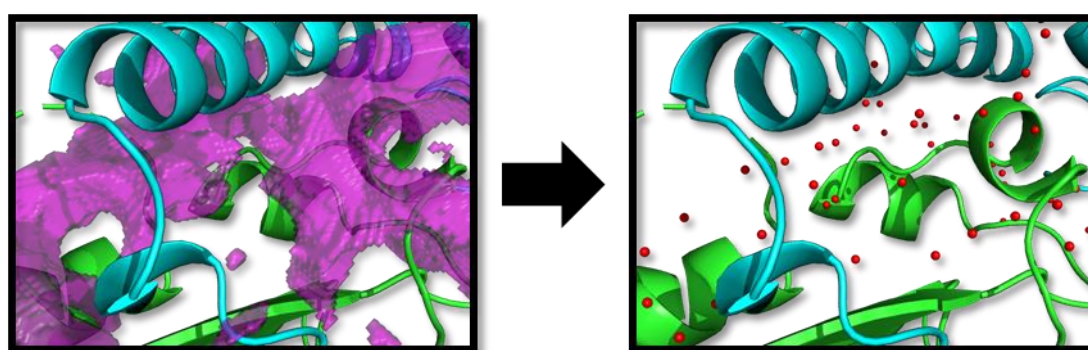


図 3. 三次元分布関数を基にした水分子の位置予測

三次元分布関数 (紫色の透過面) を基に予測水を配置することで、水分子の位置を明示的に予測することができる。図では三次元分布関数を基に配置される予測水を赤の球で示している。

水の三次元分布関数は、計算領域全体、3D-RISM 計算では全てのグリッド点について確率密度が計算されるものであるため、水分子が存在すると考えられる座標を予測するものではない。しかし、三次元分布関数から明示的な水分子の位置を予測、即ち三次元分布関数を基に水分子が存在すると考えられる座標を予測することも可能である。このような手法は三次元分布関数の有用性を更に高めるものであると期待できる。

三次元分布関数を基に、水分子の位置を明示的に予測するため、既にいくつかの手法が開発されており、代表的な手法としては確率密度が最大の位置へ水分子の配置を繰り返す「Placevent<sup>17</sup>」が挙げられる。また、近年では進化的アルゴリズムを用いた「GAsol<sup>18</sup>」などの手法も報告されている。

### 1-3-4. Placevent

ここで、Placevent による水分子の位置予測手法について説明する。Placevent は水の三次元分布関数のみを入力とし、水分子の位置を明示的に予測する手法である<sup>17</sup>。

Placevent では、始めに、三次元空間中において、最も確率密度の高い場所へ水分子が配置される。次に、予測水を配置した場所を中心に水 1 つ分の確率密度を三次元分布関数から削除し、この時点で最も確率密度の高い場所へ次の予測水を配置する。この操作を繰り返すことで、三次元分布関数を基に、水分子が存在すると考えられる位置の座標が予測される。Placevent における予測水の配置は、三次元分布関数から水分子 1 つ分の確率密度が削除できなくなるまで繰り返される（参考文献 17 Figure 1）。また、予測水が閾値（デフォルト値では 1.5）未満の確率密度の場所に配置された場合、配置を終了する。

Placevent はこのようなアルゴリズムによって水分子の位置を予測する。このため、1 つの三次元分布関数から得られるサンプルは 1 つのみであり、同一の三次元分布関数に対して Placevent を再度実行した場合でも同じサンプルが得られることとなる。Placevent による水分子の位置予測は、三次元分布関数を基に、最も尤度の高いサンプルを出力することに相当すると言える<sup>17</sup>。

### 1-3-5. 本研究の目的

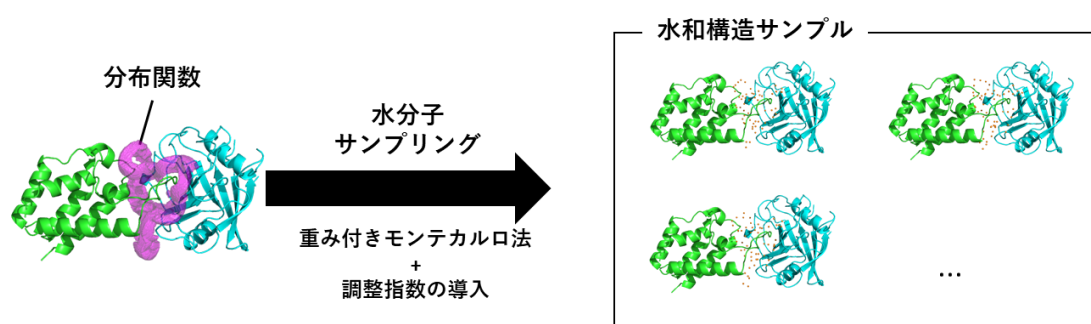


図 4. DroPred の概要

本研究では、水の三次元分布関数を基に、タンパク質周辺の水分子の位置を明示的に予測する新たな手法の開発を目指した（図 4）。



本研究室における先行研究では、重み付きモンテカルロ法を使用して水の三次元分布関数を基に水分子の位置を予測し、複数パターンのサンプルを出力する手法が開発されている<sup>19</sup>。本論文では、モンテカルロ法を用いて水分子の位置を予測したサンプルを複数出力する手法を「水分子のサンプリング」と呼ぶ。また、本論文においては以降、水分子のサンプリングにおいて水分子を配置する領域を「サンプリング領域」と呼ぶこととする。

本研究の予備調査によって、先行研究の予測精度には、改善の余地があることが示唆された。そこで本研究では、この先行研究を応用し、重み付きモンテカルロ法の重みを調整するための「調整指数」を導入した。これにより、水の三次元分布関数から得られる確率密度をより強く反映することで、より高精度に水分子のサンプリングを行う手法である「DroPred」の開発を目指した（図4）。

本研究では、3D-RISM 計算を用いて水の三次元分布関数を計算した。3D-RISMからは水の酸素原子及び水素原子の三次元分布関数が得られるが、本研究で開発した DroPred の計算では水の酸素原子の三次元分布関数を使用し、水分子中の酸素原子の座標を「水の位置」として予測した。

## 第2章. 方法

### 2-1. 本研究における重み付きモンテカルロ法

生物分子設計学教室の先行研究は、水の三次元分布関数を水分子のサンプリングへと適用したものである。

重み付きモンテカルロ法 (weighted Monte Carlo Method) <sup>20</sup> とは、モンテカルロ・サンプリングの一種であり、サンプリングにおいて、重みの大きい事象を優先的に発生させる手法である。水分子のサンプリングにおいては、重み付きモンテカルロ法の重みとして、水の確率密度を使用した。水分子のサンプリングについて考えた場合、一般的なモンテカルロ・サンプリングであれば、サンプリング領域内にランダムに予測水が配置される (図 5(a))。一方で、重み付きモンテカルロ法を用いることで、サンプリング領域の中で確率密度の高い場所を優先して予測水を配置することが可能となる (図 5(b))。

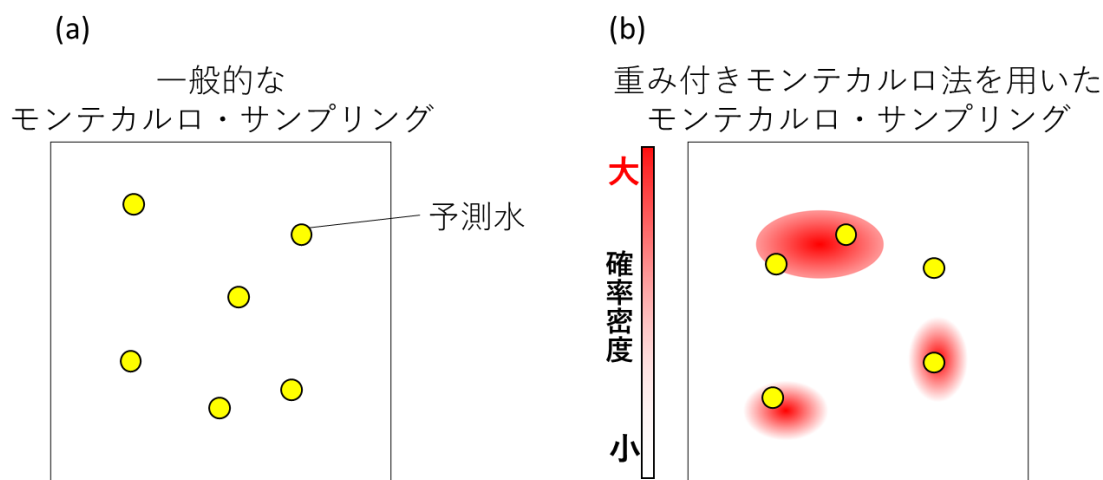


図 5. 重み付きモンテカルロ法を用いたモンテカルロ・サンプリング  
(a) は一般的なモンテカルロ・サンプリングを示す。サンプリング領域内に予測水がランダムに配置される。(b) は重み付きモンテカルロ法を用いたモンテカルロ・サンプリングを示す。予測水は確率密度の高い場所に優先的に配置される。

## 2-2. 水分子のサンプリングに対する重み付きモンテカルロ法の適用

はじめに、重み付きモンテカルロ法を、水分子のサンプリングに適用する方法について説明する。

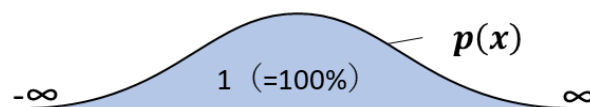


図 6. 確率密度関数の積分

ある確率変数 $x$ と確率密度関数 $p(x)$ を定義したとき、確率変数 $x$ の全値域における確率密度関数 $p(x)$ の積分値は1 (=100%) となる (図 6・式 1)。

$$\int_{-\infty}^{\infty} p(x)dx = 1 \quad (1)$$

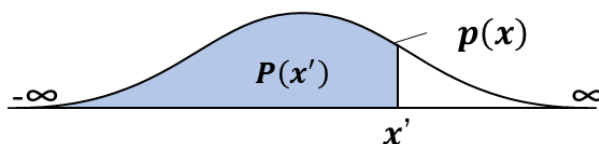


図 7. 累積確率密度

ここで、ある位置 $x'$ までの累積確率密度、即ち位置 $x'$ までの確率密度の合計値は式 2 のように表すことができる (図 7)。

$$P(x') = \int_{-\infty}^{x'} p(x)dx \quad (0 \leq P(x') \leq 1) \quad (2)$$

これによって、0 から 1 までの値と位置 $x'$ の値を対応させることができる。つまり、0 から 1 の値域で乱数を発生させることで、位置 $x'$ の値と対応させることができる (式 3)。

$$P(x') = \int_{-\infty}^{x'} p(x)dx = R \quad (3)$$

ここで、式 2 を三次元座標空間に適用し、水の三次元分布関数を確率密度関数として使用すると、以下のようなになる (式 4)。

$$P(x', y', z') = \frac{1}{N} \int_{-\infty}^{x'} \int_{-\infty}^{y'} \int_{-\infty}^{z'} \rho_0 g_0(\mathbf{r}) |_{g_0(\mathbf{r}) \geq c} dx dy dz \quad (4)$$

式 4 において  $\rho_0$  は水の酸素原子の数密度、 $g_0(\mathbf{r})$  は三次元空間上の位置  $\mathbf{r}$  における水の酸素原子の確率密度であり、 $c$  は確率密度のクライテリアである。

また、 $N$  は累積確率密度の値域を 0 から 1 にするための規格化定数であり、以下のように定義される (式 5)。

$$N = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \rho_0 g_0(\mathbf{r}) |_{g_0(\mathbf{r}) \geq c} dx dy dz \quad (5)$$

本研究室における先行研究は、式 4 を基にプログラムを構築することで、重み付きモンテカルロ法を水分子のサンプリングに適用したものである。

### 2-3. 調整重みの導入

式 4 を基に水分子のサンプリングを行った場合、確率密度が高い場所と低い場所で、予測水の配置確率に十分な差が生じないことが原因で、配置された予測水が結晶水を十分に再現しない可能性がある。そこで本研究では、式 4 における重み部分を、指数  $w$  を用いて調整した (式 6)。

$$P(x', y', z') = \frac{1}{N_w} \int_{-\infty}^{x'} \int_{-\infty}^{y'} \int_{-\infty}^{z'} \rho_0 (g_0(\mathbf{r}) |_{g_0(\mathbf{r}) \geq c})^w dx dy dz \quad (6)$$

これにより、確率密度が高い場所にはより大きな重みが、確率密度が低い場所へはより小さな重みが与えられることとなり、確率密度の高い場所をより優先して予測水を配置することが可能となる。本研究では、重みの調整を行うための指数  $w$  を「調整指数」、調整指数によって調整された重みを「調整重み」とそれぞれ定義した。

規格化定数 $N_w$ は、次のように再定義される（式7）。

$$N_w = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \rho_0 (g_0(\mathbf{r}) |_{g_0(\mathbf{r}) \geq c})^w dx dy dz \quad (7)$$

本研究では、水の三次元分布関数を計算するために、前述した 3D-RISM 計算を使用した。3D-RISM 計算における三次元分布関数は、計算領域内の各グリッド点 $(x_i, y_j, z_k)$ に確率密度 $g_0(x_i, y_j, z_k)$ として与えられる。本研究では、3D-RISM 計算によって与えられた確率密度を基に、確率密度 $G_0(x_i, y_j, z_k)$ を再定義した（式8）。

$$G_0(x_i, y_j, z_k) = \frac{1}{8} \sum_{l=0}^1 \sum_{m=0}^1 \sum_{n=0}^1 (g_0(x_{i+l}, y_{j+m}, z_{k+n}) |_{g_0(x_{i+l}, y_{j+m}, z_{k+n}) \geq c})^w \quad (8)$$

これにより、式6、式7は以下のように再定義される（式9・式10）。

$$P(X', Y', Z') = \frac{1}{N_w} \sum_{i=0}^{X'} \sum_{j=0}^{Y'} \sum_{k=0}^{Z'} \rho_0 V G_0(x_i, y_j, z_k) \quad (9)$$

$$N_w = \sum_{i=0}^X \sum_{j=0}^Y \sum_{k=0}^Z \rho_0 V G_0(x_i, y_j, z_k) \quad (10)$$

ここで、 $V$ は 3D-RISM 計算が実行されたボックスを分割するボクセルの体積、 $X$ 、 $Y$ 、 $Z$ はそれぞれ  $X$  軸、 $Y$  軸、 $Z$  軸方向のグリッド数である。

本研究では、点 $(0,0,0)$ から点 $(X', Y', Z')$ までの累積確率である式9をベースにプログラムを構築し、水分子のサンプリングを行うプログラムである「DroPred」を開発した。

## 2-4. プログラム構築

前述した式 9 を基に、水分子のサンプリングを行うプログラムを開発した。本研究におけるプログラムの開発は、処理速度の速いコンパイル言語である C 言語を使用して行った。

### 2-4-1. 予測水同士の距離

予測手法によって予測水を配置する際に、配置される他の予測水の影響を考慮しない場合、2 つ以上の予測水を非常に近い位置に配置することも可能である。しかし、水分子も物理的な大きさを持っているため、実在するタンパク質中では 2 つ以上の水分子が極端に近い位置に存在することは考えづらい。このため、水分子の位置を予測する際には、予測水同士が過剰に近い位置に配置されないようにするアルゴリズムを導入することが好ましいと考えられる。Placevent では、予測水を配置した際に、三次元分布関数から水分子 1 つ分の確率密度が取り除かれるが、これが過度に近い位置に予測水が配置されないようなアルゴリズムにも相当すると述べられている<sup>17</sup>。DroPred においては、ある位置に予測水を配置する際、その位置と既に予測水が配置された全ての位置の距離を計算し、既にその周辺に予測水が配置されているのであれば予測水を配置しない、というアルゴリズムを組み込むことで、予測水同士が極端に近い位置に配置されることを回避できるようにした。

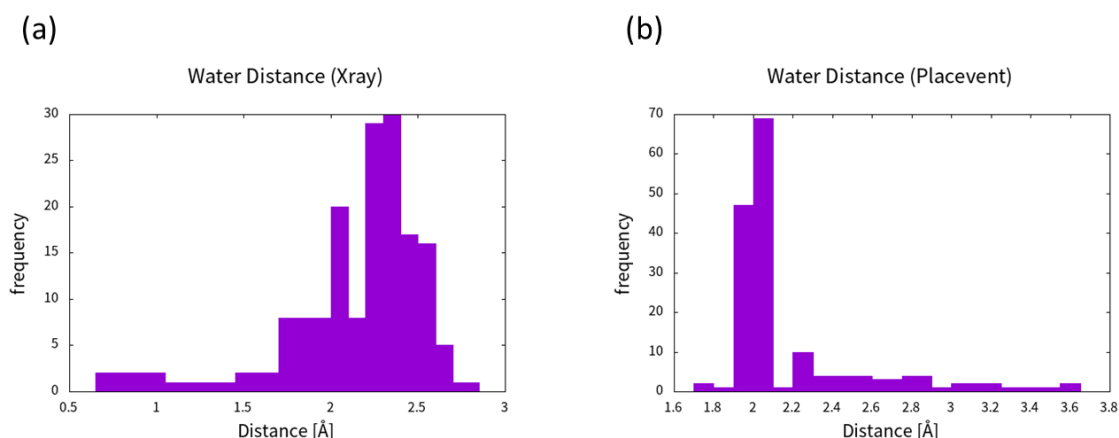


図 8. 水分子同士の最小距離

水分子同士の最小距離のヒストグラム。(a)は X 線結晶構造解析によって得られた実験構造中の結晶水同士の最小距離、(b)は Placevent によって予測された水分子同士の最小距離である。横軸は最小距離 (Distance) を、縦軸は頻度 (frequency) をそれぞれ示している。

ここで、このアルゴリズム中において、「ある位置の周辺に予測水が既に配置されている」という状態を定義した。そのために、X線結晶構造解析によって解析された B-factor が 40 未満の水分子、及び Placevent による予測で得られるサンプルに含まれる予測水について、酸素原子同士の距離を総当たりで計算し、最小距離を調査した (図 8 (a))。調査は本研究の性能評価に使用したテストセット (151 ターゲット) を使用して行った。テストセットについては第 2-5-2 章で説明する。この結果、X線結晶構造解析による構造では、結晶水同士の最小距離は 2.0 から 2.5 Å、Placevent による予測水同士の間の最小距離は 2.0 Å 前後の頻度が高いことが分かった。この結果を基に、本研究においては 2.0 Å 以内に既に予測水が配置されているグリッド点を、「予測水が既に配置されているグリッド点」として定義した。

## 2-4-2. DroPred のアルゴリズム

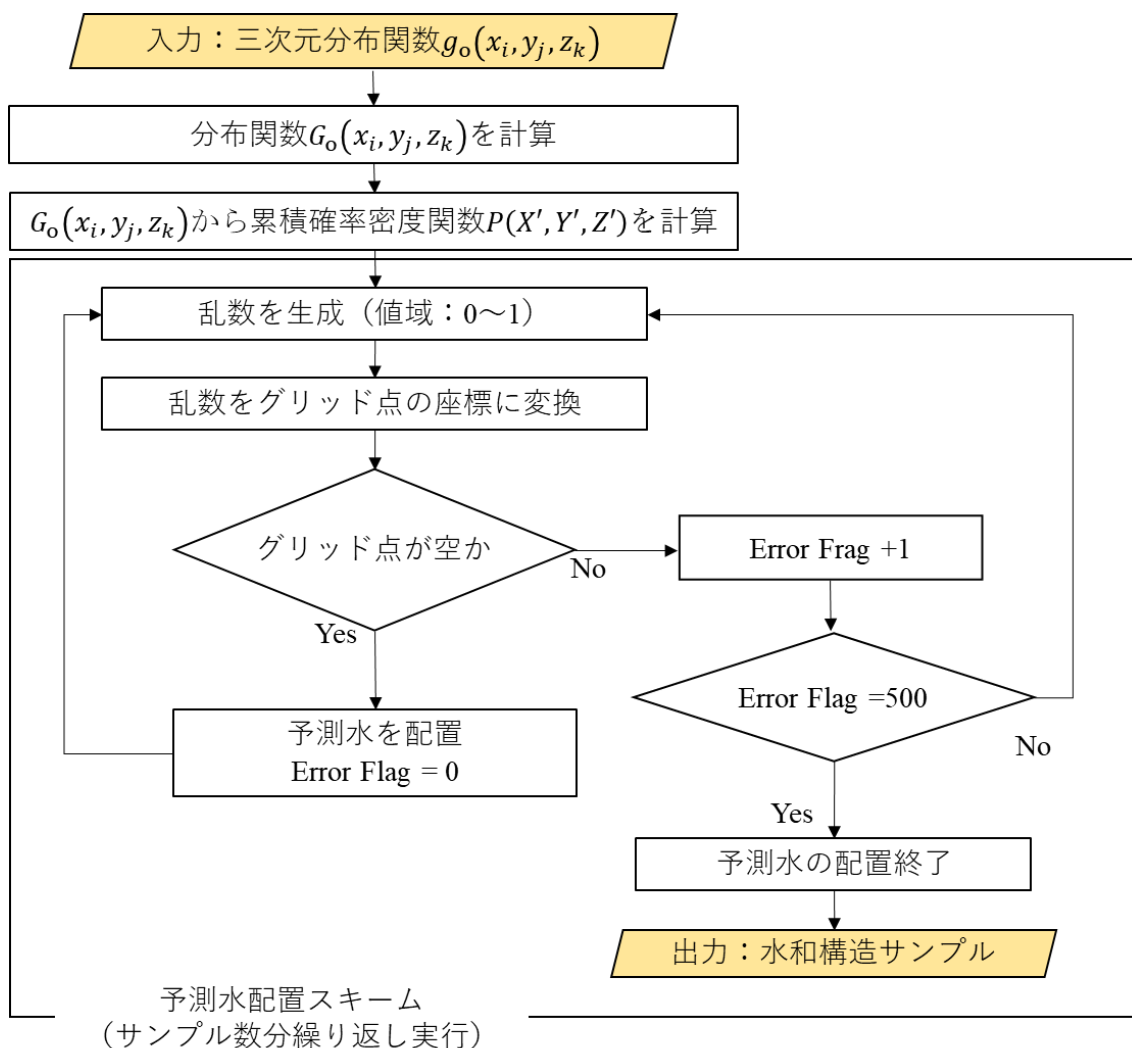


図 9. DroPred による水分子のサンプリングアルゴリズム

DroPred は、図に示したアルゴリズムによって予測水を配置することで、水分子のサンプリングを実行する。

DroPred によって水分子のサンプリングを行うためのアルゴリズムについて説明する。(図 9) DroPred は、水の酸素原子の三次元分布関数のみを入力として、水分子のサンプリングを行うことができる。

入力として水の三次元分布関数が与えられると、与えられた三次元分布関数は累積確率密度関数に変換され、予測水の配置スキームが実行される。本研究においては、確率密度のクライテリアを、 $g_o(\mathbf{r}) \geq 2$ とした。予測水の配置スキーム



ムでは、初めに 0 から 1 の値域で乱数が生成され、前述した式 9 に基づいて乱数はボックス内のグリッド点の座標へと変換される。次に、グリッド点と、既に配置された全ての予測水との距離が計算され、周囲 2.0 Å 以内に既に配置された予測水が存在しなければ、そのグリッド点は空であると定義される。グリッド点が空である場合、そのグリッド点に予測水を配置し、再び乱数の生成からスキームを実行する。グリッド点が空でなかった場合は、予測水の配置は行わず、同様に再び乱数の生成からスキームを実行する。配置された水分子の数が多くなるほど、予測水の配置は失敗しやすくなる。予測水の配置が 500 回連続で行われなかった場合、既に十分な数の予測水が配置されたとして、サンプリングを終了する。

予測水の配置スキームは、乱数のシード値を変更することで、実行するたびに毎回異なるサンプルを得ることができる。すなわち DroPred では、この予測水の配置スキームを任意の回数だけ繰り返すことで、任意の数の異なるサンプルを出力することが可能である。

### 2-4-3. 並列化による DroPred 計算の効率化

前述した通り、DroPred は予測水の配置スキームを繰り返すことで、任意の数のサンプルを出力することが可能である。しかしながら、多数のサンプルの作成を行う場合は計算コストが大きくなる。そこで、共有メモリ型マシンで並列プログラミングを可能にする API である「OpenMP」を用いた並列処理を DroPred プログラム中に実装することで、予測水の配置スキームを並列化して実行可能にすることで、計算時間を短縮した。

## 2-5. 性能評価

本研究において開発した DroPred をタンパク質 - タンパク質複合体界面に適用し、性能評価を行った。

### 2-5-1. 性能評価プロトコル

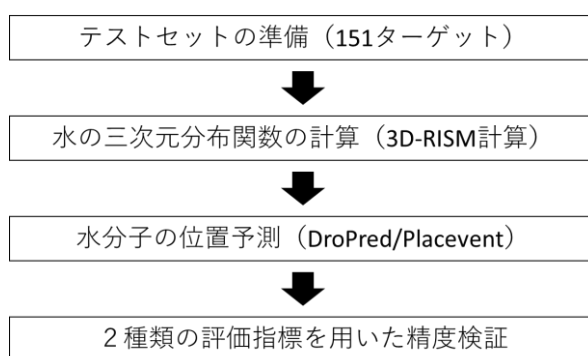


図 10. DroPred の性能評価プロトコル

DroPred の性能評価を行うためのプロトコルについて説明する。(図 10)

初めに、性能評価のためのタンパク質 - タンパク質複合体のテストセット (151 ターゲット) を準備した。テストセットの各ターゲットについて、3D-RISM 計算を実行し、水の三次元分布関数を得た。得られた水の三次元分布関数に対して、DroPred 及び比較検証のための Placevent を実行し、水分子の位置予測を行った。最後に、2 種類の評価指標を用いて DroPred から得られるサンプルの精度検証を実行した。

## 2-5-2. テストセットの準備

創薬研究において、タンパク質-タンパク質複合体 (PPI) は重要なターゲットとなっている。そこで本研究では、DroPred の性能評価をタンパク質-タンパク質複合体界面に存在する水分子によって行うこととした。

DroPred の性能評価を行うために、151 個のターゲットからなるタンパク質 - タンパク質複合体のテストセットを準備した<sup>21,22</sup>。準備したテストセットは、X線結晶構造解析によって得られた 2 量体 (103 ターゲット)、3 量体 (47 ターゲット)、5 量体 (1 ターゲット) のタンパク質-タンパク質複合体から構成されており、全てのターゲットの分解能は 2.5 Å より良好であった。複合体の残基数は 152 から 809 であり、界面残基数は 67 から 413 であった。準備したテストセットに含まれる全てのターゲットの PDB ID については、Supporting Information 中の表 S1 に記載した。

性能評価に使用する結晶水は、タンパク質 - タンパク質複合体の界面に存在する水分子とした。タンパク質 - タンパク質複合体界面は、「2 つ以上のポリペプチド鎖からそれぞれ 5.0 Å 以内の距離に存在する領域」と定義した。また、分子の熱振動の大きさを表す値である B-factor<sup>23,24</sup> が 40 以上の結晶水については、タンパク質 - タンパク質複合体の形成に関与しない結晶水である可能性を考慮し、本研究における性能評価においては使用しなかった。

テストセットの各ターゲットには、7 個から 119 個、平均で 36 個の結晶水が含まれていた。また、テストセット全体では、5448 個の結晶水を性能評価の対象とした。

尚、本研究において使用したタンパク質の立体構造、及び結晶水の座標データは、全て PDB から取得したものであり、全てのタンパク質及び結晶水は同一の PDB ファイル内に記載されているものを使用している。

## 2-5-3. 水の三次元分布関数の計算 (3D-RISM 計算)

テストセットの 151 ターゲットそれぞれに対して、3D-RISM 計算を実行し、水の三次元分布関数の計算を行った。

3D-RISM 計算の入力として、前述したテストセットの PDB ファイルから、水素原子を除くタンパク質の座標データのみを抽出したファイルを作成した。これ

を入力とし、AmberTools<sup>25</sup>の「sander」プログラムを使用して3D-RISM計算を実行した。本研究ではClosure方程式としてKovalenko-Hirata (KH) Closure方程式を使用した。第1章で述べたように、3D-RISM計算は、ターゲットを三次元グリッド空間に配置し、タンパク質を中心とした十分な大きさの直方体のボックスを計算領域として行い、ボックス内の全てのグリッド点について確率密度が計算される。三次元グリッド空間のグリッド幅をより小さくすることで、水の確率密度をより詳細に計算することが可能になるが、その分計算コストは大きくなる。本研究では、三次元グリッド空間のグリッド幅を0.25 Åとした。また、水素原子を除くタンパク質分子からボックスの縁、即ちバルクまでの最短距離(バッファ)は20 Åとした。力場として、水にはcSPCE、タンパク質にはff99SB<sup>26</sup>を使用し、水のモル濃度は55.5 Mとした。収束計算は、収束条件を10e-5とし、最大10,000回の反復計算を実行した。

尚、3D-RISM計算における計算コストは、計算を行うグリッド点の個数、つまり計算領域の体積に強く依存する。このため、3D-RISMの計算時に用意するボックスの体積が小さくなるほど、3D-RISM計算の計算コストは小さくなる。このボックスは、各辺が三次元グリッド空間のXYZ軸と並行となるように設定されるため、ターゲットの向きによってはボックスの大きさに無駄が生じる可能性がある。そこで本研究では、3D-RISM計算を行う前に必要なボックスの体積が最小となるようターゲットの向きと位置を調整し、3D-RISM計算を効率よく実行できるようにした。具体的には、①ターゲットのタンパク質内で最も距離の離れた2原子を結ぶ直線が三次元座標空間上のZ軸と一致するように構造を移動し、②Z軸を中心に、XY平面の面積が最小となるように構造を回転させる、といった2段階の操作によって構造の向きを最適化した。これにより、例えば、あるターゲット(PDB ID: 1JTD)において、PDBに登録された構造に対して3D-RISM計算を実行するためには、98 (Å) × 112 (Å) × 105 (Å) = 1152480 (Å<sup>3</sup>)のボックスが必要であったが、調整後に必要なボックスの大きさは98 (Å) × 90 (Å) × 120 (Å) = 1058400 (Å<sup>3</sup>)であり、ボックスの体積を約8%減少させることに成功した。

これにより、テストセットの各ターゲットについて計算された水の三次元分布関数を得た。

#### 2-5-4. 水分子の位置予測

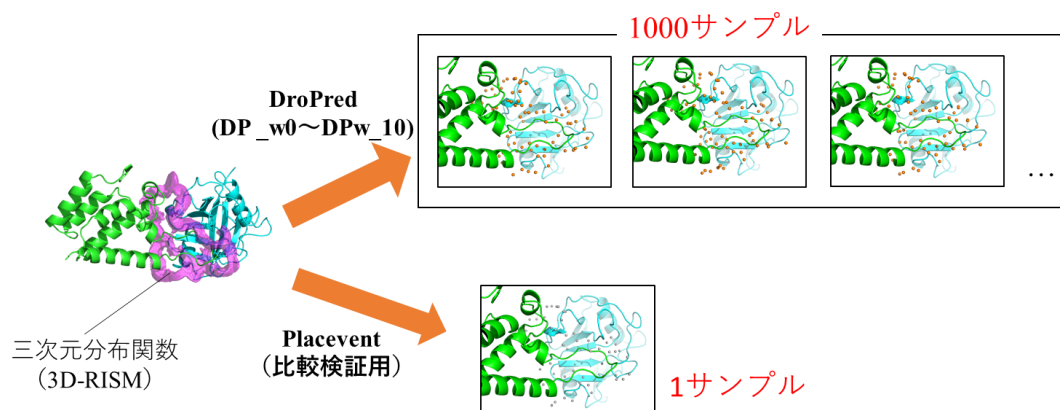


図 11. DroPred・Placevent による水分子の位置予測の実行

3D-RISM 計算によって得られた水の三次元分布関数を基に、DroPred・Placevent による水分子の位置予測を実行した。

テストセットの 151 ターゲットに対して、3D-RISM 計算によって得られた水の三次元分布関数を基に、水分子の位置予測を実行した (図 11)。

本研究においては、予測水の配置を行うサンプリング領域をタンパク質 - タンパク質複合体界面とした。テストセットと同様、タンパク質 - タンパク質複合体界面は、「2 つ以上のポリペプチド鎖からそれぞれ 5.0 Å 以下の距離に存在する領域」と定義した。サンプリング領域を限定するため、3D-RISM 計算によって得られた水の三次元分布関数に対して前処理を行い、界面領域以外のグリッド点の確率密度を 0 とした。

第 2 章にて述べたように、DroPred は分布関数から得られる確率密度を基にした重みを、調整指数を用いて調整することで、重みの調整を行わない場合と比較して、より確率密度の高い場所を優先して予測水を配置することが可能である。本研究では、調整指数ごとに DroPred の性能を評価することも目的の 1 つとした。この目的のため、0 から 10 までの整数値、11 種類の調整指数を用いて、DroPred による水分子の位置予測を実行した。本論文においては、用いた調整指数ごとに、DP\_w0 から DP\_w10 とそれぞれ呼称することとする。ここで、調整指数 0 を用いた DP\_w0 は、一般的なモンテカルロ法によって、サンプリング領域内に予測水をランダムに配置することに相当する。また、調整指数 1 を用いた DP\_w1 は、三次元分布関数から得られる確率密度をそのまま重みとして使用した重み付きモンテカルロ法によるサンプリングであり、先行研究において開発された手法に相当する<sup>19</sup>。調整指数 2 から 10 を用いた DP\_w2 から DP\_w10 は、重み

を調整した重み付きモンテカルロ法によるサンプリングであり、大きな調整指数を用いるほど、予測水が配置されるスキーム（図 9）において、確率密度の高い場所により高い確率で予測水が配置されるようになる。

また、DroPred は、任意の数のサンプルを出力することが可能であるが、本研究では、各ターゲットについて DP\_w0 から DP\_w10 の手法でそれぞれ 1000 サンプルずつを出力した。DroPred の出力サンプル数については、第 3 章において追加で説明する。

DroPred による結晶水の再現度について比較検証を行うため、同一の三次元分布関数を使用して、代表的な予測手法である Placevent による水分子の位置予測も実行した。Placevent の実行は、デフォルトのパラメータを使用して行った。Placevent からは、各ターゲットに対して 1 つのサンプルのみが出力される。

#### 2-5-5. DroPred の性能評価

DroPred (DP\_w0 から DP\_w10)、及び Placevent によって水分子の位置予測を行ったサンプルを用いて、DroPred の性能評価を行った。DroPred の性能は、「結晶水の位置の再現度」及び「分布関数の再現度」の 2 種類の評価指標によって評価した。

評価方法及び評価指標については、第 3 章において詳しく説明する。

## 第3章. 結果・考察

この章では、第2章で開発を行った DroPred の性能評価の結果を示す。

### 3-1. 水分子の再現の定義

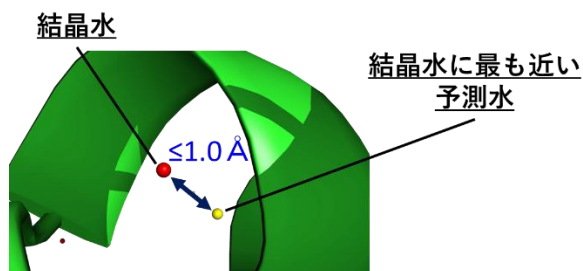


図 12. 結晶水の再現の定義

赤の球は結晶水を、黄色の球はこの結晶水に最も近い場所に配置された予測水をそれぞれ示している。結晶水が存在する場所の周囲  $1.0\text{\AA}$  以下の場所に予測水を配置することが出来た場合、その結晶水は「再現された」と定義した。

結晶水の位置の再現度を評価するため、「結晶水の再現」を定義した。図 12 では、結晶水が存在する場所を赤色の球で、配置した予測水のうち、最も結晶水に近い位置に配置されたものを黄色の球でそれぞれ示している。図に示すように、あるサンプルにおいて、結晶水が存在する場所の周囲  $1.0\text{\AA}$  以下の場所に予測水を配置することが出来た場合、その結晶水は「再現された」と定義した。

### 3-2. DroPred の出力サンプル数

DroPred のアルゴリズムでも述べたように、DroPred からは、任意の数のサンプルを出力することが可能である。本研究では各ターゲットにつき、1000 サンプルずつを出力した。

表 1: 1 つ以上のサンプルによって再現された水分子の割合の推移

サンプル数	1 つ以上のサンプルによって再現された結晶水の割合(%)										
	DP_w0	DP_w1	DP_w2	DP_w3	DP_w4	DP_w5	DP_w6	DP_w7	DP_w8	DP_w9	DP_w10
100	99.02	98.76	98.53	97.94	97.33	96.31	94.52	93.74	90.16	85.85	79.56
200	99.41	98.99	98.82	98.42	98.24	97.41	96.37	95.37	92.69	89.22	83.39
300	99.46	99.08	98.88	98.69	98.53	97.88	97.08	96.16	94.09	90.65	85.48
400	99.46	99.08	98.93	98.76	98.72	98.13	97.42	96.75	94.71	91.66	86.77
500	99.49	99.10	98.97	98.81	98.80	98.35	97.57	96.94	95.10	92.36	87.66
600	99.49	99.14	99.01	98.90	98.85	98.48	97.83	97.16	95.49	92.96	88.34
700	99.51	99.17	99.05	98.92	98.87	98.60	97.88	97.34	95.70	93.26	88.84
800	99.51	99.21	99.05	98.98	98.89	98.63	98.02	97.44	96.01	93.55	89.22
900	99.51	99.22	99.05	99.04	98.91	98.73	98.14	97.56	96.15	93.84	89.51
1000	99.51	99.22	99.08	99.04	98.93	98.76	98.18	97.64	96.22	94.15	89.84

DroPred によってあるサンプル数を出力した場合に、性能評価に使用した 5448 個の結晶水のうち、1 つ以上のサンプルで再現された予測水の割合の推移。サンプル数が 100 増加した場合に、1 サンプル以上で再現された結晶水の個数が増加した割合が 1%未満であったものをオレンジ色で着色した。

表 1 は、DroPred によってあるサンプル数を出力した場合に、性能評価に使用した 5448 個の結晶水のうち 1 つ以上のサンプルによって再現された結晶水の個数の割合の推移を示したものである。出力サンプル数が 100 増加した場合に、1 サンプル以上で再現された結晶水の割合が 1%未満であったものを、オレンジ色で着色して示した。表 1 に示すように、1000 サンプルを出力した時点で、DP\_w0 から DP\_w9 において、評価に使用した結晶水の 90%以上が 1 サンプル以上で再現されており、かつ DP\_w0 から DP\_w10 までの全てにおいて、新たに再現されるようになる結晶水の割合が 1%未満となった。



表 1 は、結晶水が存在する 5448 か所の位置において、サンプリング中に一度でも予測水が配置されるのかを示した結果とも考えることができる。どの手法においても、出力サンプル数を増加させると、再現できる結晶水の位置が増加するが、限界が存在し、1000 サンプルを出力しても数か所は再現することのできない結晶水の位置が存在することが示された。

### 3-3. 予測水の配置数

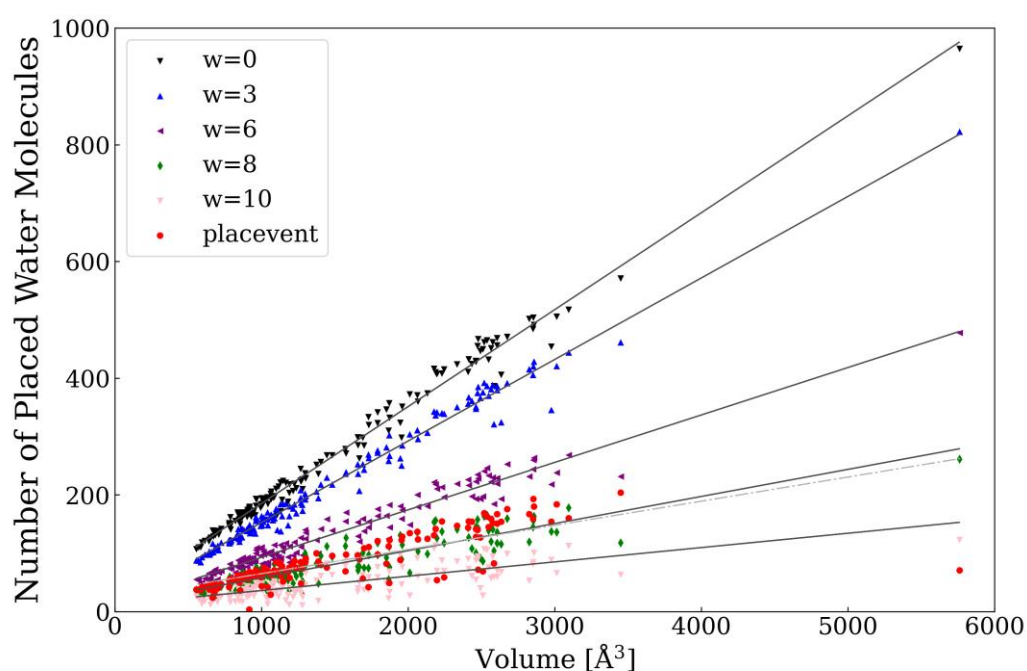


図 13. 各手法における予測水の配置数

各手法における予測水の配置数。横軸はターゲットのタンパク質-タンパク質界面において 0 より大きな確率密度を持つ領域の体積 (Volume) を、縦軸は予測水の配置数 (Number of Placed Water Molecules) を示す。

DroPred のアルゴリズムでも述べたように、DroPred の予測水配置スキームにおいては、終了条件を満たすまで予測水が配置される。このため、同じターゲットに対するサンプルであっても、配置される予測水の個数は異なる。

図 13 に、DroPred 及び Placevent による予測水の配置数を示す。DroPred については、151 ターゲットについて、1000 サンプル分の配置数の平均値をプロットしている。Placevent からは各ターゲットについて 1 サンプルのみが出力さ

れるため、各ターゲットの配置数をそのままプロットしている。DroPred は、使用する調整指数が大きくなるほど配置数が減少し、DP\_w8 における配置数は Placevent の配置数とほぼ同等であった。DP\_w0 から DP\_w7 の配置数は Placevent よりも多く、DP\_w9、DP\_w10 の配置数は Placevent と比較して少数であった。各ターゲットに対する水の配置数に関する詳細なデータについては、Supporting Information 中の表 S2 に記載した。

大きな調整指数を使用した場合に配置される予測水の数が減少した原因としては、確率密度の低い場所に予測水が配置される頻度が減少したことで、調整指数が小さい場合と比較して、終了条件が満たされるまでに確率密度の低い場所へ配置される予測水の数が減少したことが大きな原因であると考えられる。

また、結晶水の再現度を評価する場合、基本的には多くの予測水が配置されているサンプルほど、多くの結晶水を再現している可能性が高くなると考えられる。しかし、予測水を多く含むサンプルには、結晶水を再現していない予測水も多く含まれている可能性も高い。そこで本研究では、DroPred の性能を Placevent と比較するために、各ターゲットに対して Placevent と同数までの予測水を、評価に使用することとした。DroPred による各サンプルに Placevent を超える数の予測水が含まれていた場合、配置が行われた順に Placevent と同数までを評価に使用した。尚、サンプルに含まれる予測水が Placevent による予測水の個数を下回っている場合については、全ての予測水の評価に使用した。

### 3-4. 評価 1： 結晶水の位置の再現度による評価

#### 3-4-1. 評価の概要

はじめに、DroPred が実験によって決定された結晶水の位置を精度よく再現できているかを評価した。

タンパク質周辺の水分子の予測手法を開発する目的の 1 つとして、実験的な手法の代替として使用し、予測された水分子の情報を使用して、タンパク質の機能や相互作用などの解析を行うことが挙げられる。実験的な手法で得られるタンパク質周辺の水分子を精度よく再現できる予測手法であれば、実験的な手法が適用できないタンパク質や、予測構造などの仮想的なタンパク質に対しても適用可能であると考えられる。即ち、予測手法によって得られるサンプルは、実験的な手法で得られる水分子の位置を精度よく再現できることが好ましい。こ

のため、DroPredによって予測された水分子が実験的な手法で得られた結晶水をどの程度再現できるかを、DroPredの性能評価における基準の1つとした。

DroPredはサンプリング手法であるため、様々なパターンのサンプルが得られることで、出力した1000サンプルの中で、結晶水の再現度に幅が生じることが予測される。本研究では主に、サンプル全体としての結晶水の再現度を考えるために1000サンプルの中央値に着目した。また、出力したサンプルの中に再現度の高いサンプルが含まれていれば、サンプルのスコアリングなどの手法を使用することで、再現度の高いサンプルを選び出すことが出来る可能性があることから、出力したサンプルの中で最も予測水の再現度が高いサンプルにも着目した。

また、結晶水の再現度の基準として、Placeventによる結晶水の再現度を比較対象とした。1-3-4章で述べた通り、Placeventは水の三次元分布関数から水分子の位置を予測する代表的な手法であり、三次元分布関数のみを使用して水分子を配置する手法である。また、確率密度が高い場所から予測水を配置するというアルゴリズムは、確率密度が高い場所を優先して予測水を配置するDroPredのアルゴリズムと類似した部分があると言える。このため、Placeventと同等以上の結果が得られるのであれば、DroPredによるサンプルは結晶水を十分に再現できると考えられる。

### 3-4-2. Coverage の定義

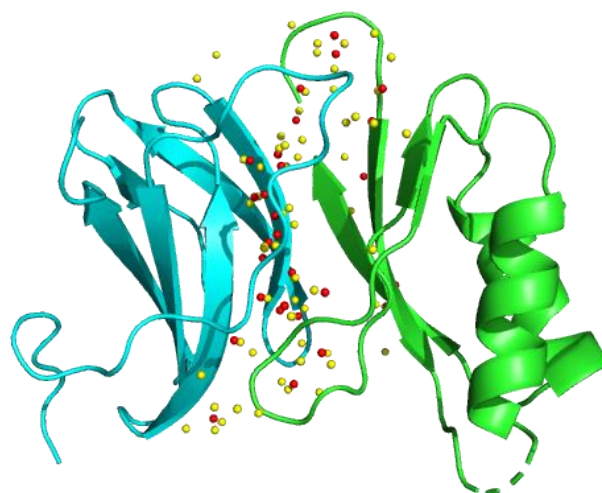


図 14. 結晶水と予測水の位置の比較

タンパク質-タンパク質界面 (PDB ID: 5A6W) における水分子の位置を DP\_w6 で予測したサンプルの 1 つ (サンプル番号: 0)。赤の球で結晶水を、黄色の球で DP\_w6 による予測水を示した。

あるサンプルの結晶水の再現度を表す評価指標として、「coverage」を定義した<sup>21,27</sup>。Coverage は、あるターゲットに含まれる結晶水のうち、そのサンプルによって再現されたものの割合として定義した。例として、図 14 にターゲットの 1 つ (PDB ID: 5A6W) を示す。このターゲットには、結晶水 (赤球) が 28 個含まれていたが、DP\_w6 によって作成されたサンプル (サンプル番号 0) の予測水 (黄球) によって、16 個の結晶水が再現された。このため、このサンプルの coverage は  $16/28 = 57.1\%$  となる。

### 3-4-3. ターゲット毎の結果

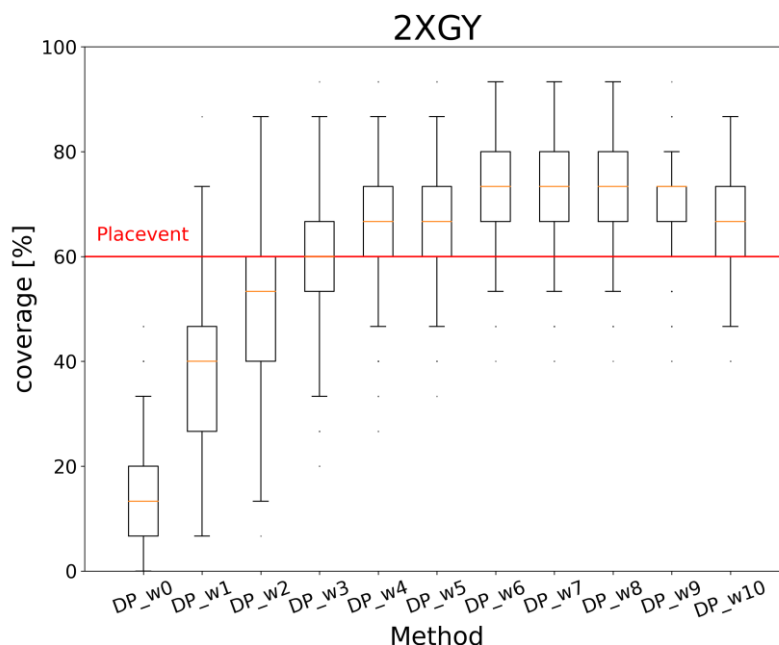


図 15. 2XGY におけるサンプリング結果

サンプリングの結果 (PDB ID: 2XGY)。横軸は DP\_w0 から DP\_w10 までの各予測手法 (Method) を、縦軸は coverage を示す。1000 サンプルの coverage の中央値をオレンジ色で示した。比較のため、2XGY に対する Placevent のサンプルの coverage を赤の線で示した。

各ターゲットに対する結晶水の再現度による評価結果をボックスプロットで示した。全てのターゲットのデータは、Supporting Information 中の図 S1 に記載した。ここでは、その中から抜粋して結果を示す。

はじめに、DroPred による予測が良好な結果を示したターゲット (PDB ID: 2XGY) についての結果を示す (図 15)。

このターゲットについて Placevent で水の位置予測を行ったサンプルの coverage は 60.0%であった。DP\_w1 では、DP\_w0 と比較すると、coverage が大きく上昇していることが確認できるものの、1000 サンプルの coverage の中央値は 40.0%と低く、結晶水が十分に再現出来ていないとは言えない。ここで、DP\_w2 以降の結果を見ると、DP\_w1 と比較して、coverage が更に上昇し、coverage の中央値は DP\_w6 で最大値 (73.3%) となった。DP\_w4 以降では、coverage の中央値

が Placevent を上回っており、ボックスプロットの上端を見ると、1000 サンプルの中には coverage が Placevent を大きく上回るサンプルも含まれていることが分かる。

この結果は、調整指数を用いた重みの調整によってサンプリングにおける結晶水の位置の再現度が向上することを示唆していると考えられる。

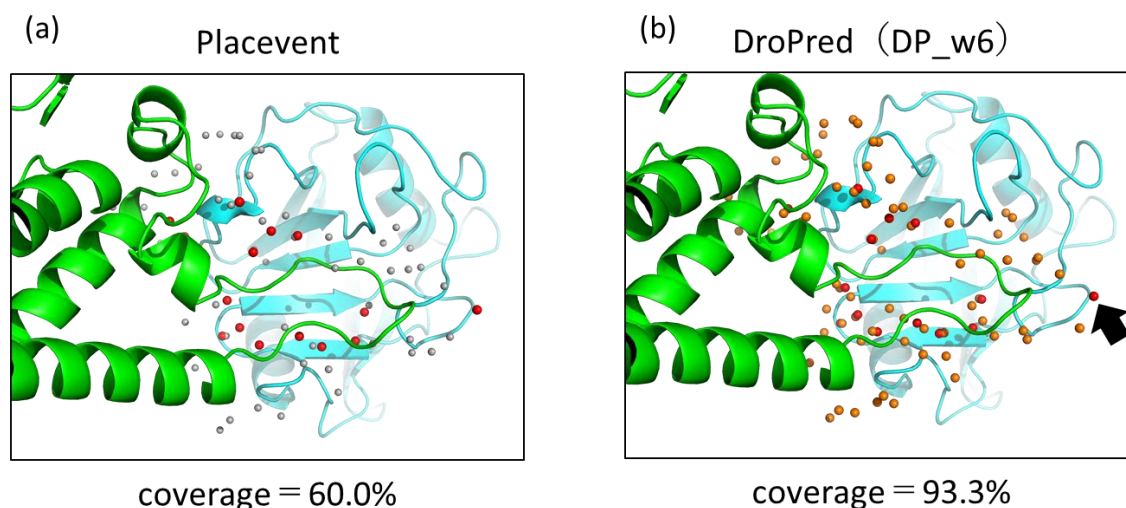


図 16. サンプルの立体構造

タンパク質-タンパク質界面(PDB ID:2XGY)の水分子の位置を Placevent・DroPred の各手法で予測した結果。(a) Placevent によるサンプル。結晶水を赤色の球で、Placevent による予測水を白色の球でそれぞれ示した。(b) DroPred によって出力されたサンプルのうち、最も coverage が良好であったサンプル(DP\_w6)。DroPred による予測水をオレンジの球で示した。

DroPred による水分子のサンプリングが良好な結果を示したターゲット (PDB ID: 2XGY) について、実際の立体構造を示す (図 16)。このターゲットには、赤色の球で示すように、結晶水が 15 個含まれていた。図 16(a)は、このターゲットに対する Placevent による水分子の位置予測結果である。Placevent による水分子の位置予測では、予測水が 57 個配置された。これによって 15 個中の 9 個の結晶水が再現され、coverage は 60.0%であった。一方、DroPred によって出力されたサンプルのうち、最も coverage が良好であったサンプル (DP\_w6) を図 16(b)に示す。このサンプルにおいては、図中に矢印で示した結晶水以外、全ての結晶水が再現され、coverage は 93.3%であった。このように、DroPred から得られるサンプル中には、Placevent の coverage を大きく上回るサンプルが含まれる可能性があることが示唆された。

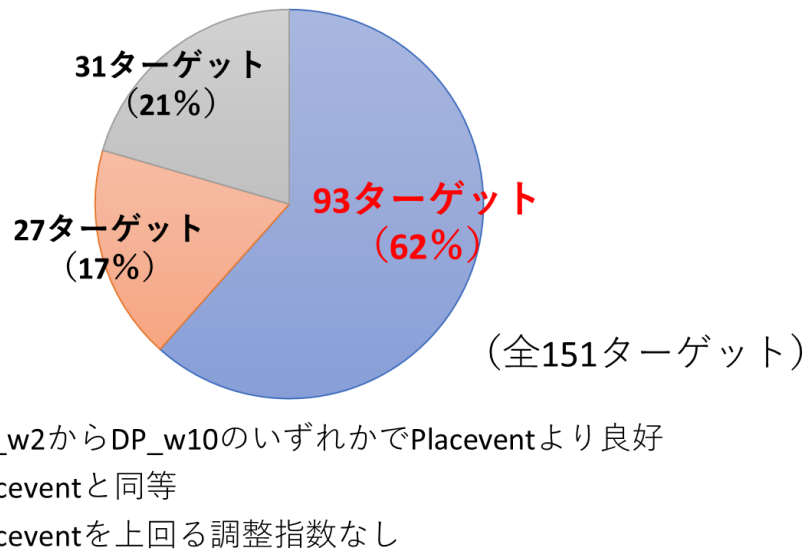


図 17. DroPred と Placevent の coverage の中央値比較

DroPredによって出力したサンプルの coverage の中央値を Placevent の coverage と比較し、coverage の中央値が Placevent を上回る調整指数が存在するか検証した。

テストセットの 151 ターゲットについて、DroPred (DP\_w2-DP\_w10) によって出力した 1000 サンプルの coverage の中央値と Placevent によるサンプルの coverage を比較し、DroPred による 1000 サンプルの coverage の中央値が Placevent によるサンプルの coverage を上回るような調整指数が存在するかを調査した (図 17)。その結果、全体の 62%にあたる 93 ターゲットで、coverage の中央値が Placevent を上回る調整指数が存在した。17%にあたる 27 ターゲットでは、coverage の中央値の最大値が Placevent と同等であり、全体の 79%のターゲットにおいて、coverage の中央値が Placevent と同等以上となる調整指数が存在した。一方で、21%にあたる 31 ターゲットにおいては、いずれの調整指数を用いた場合でも、coverage の中央値が Placevent を上回らなかった。

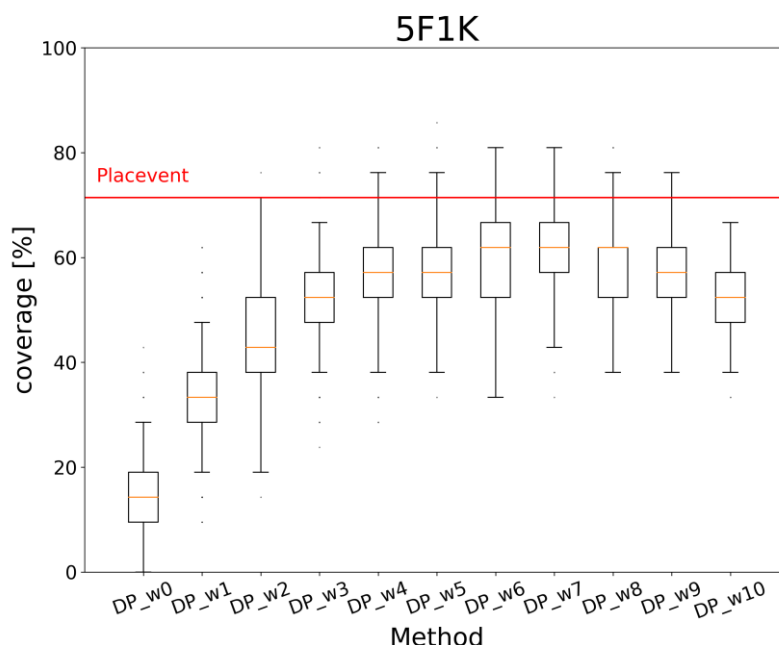


図 18. 5F1K におけるサンプリング結果

サンプリングの結果 (PDB ID : 5F1K)。横軸は DP\_w0 から DP\_w10 までの各予測手法 (Method) を、縦軸は coverage を示す。1000 サンプルの coverage の中央値をオレンジ色で示した。比較のため、5F1K に対する Placevent のサンプルの coverage を赤の線で示した。

DroPred において、どの調整指数を用いた場合でも coverage の中央値が Placevent を上回らなかったターゲットの 1 つ (PDB ID : 5F1K) についての結果を示す (図 18)。このターゲットに対する Placevent によるサンプルの coverage は 71.4%であった。DroPred による 1000 サンプルの coverage の中央値は DP\_w6 から DP\_w8 で最大 61.9%であり、Placevent のサンプルの coverage を下回った。しかし、DP\_w2 以降では、DP\_w1 に比べて coverage の大幅な上昇が見られた。また、DP\_w6、DP\_w7 などでは、ボックスプロットの上端が Placevent の coverage (71.4%) を超えており、出力した 1000 サンプル中には、coverage が Placevent を超えるサンプルが含まれていたことが分かる。重みの調整による DP\_w2 以降における coverage が上昇すること、及び出力したサンプルの中に Placevent の coverage を超えるサンプルが含まれることは、coverage の中央値が Placevent を上回らなかったターゲットを含む全てのターゲットにおいて確認することができた。この結果は、ターゲットの種類に関わらず、調整指数を用いた重みの調整によって、一定の効果が得られることを示唆していると考えられる。



### 3-4-4. テストセット全体の傾向

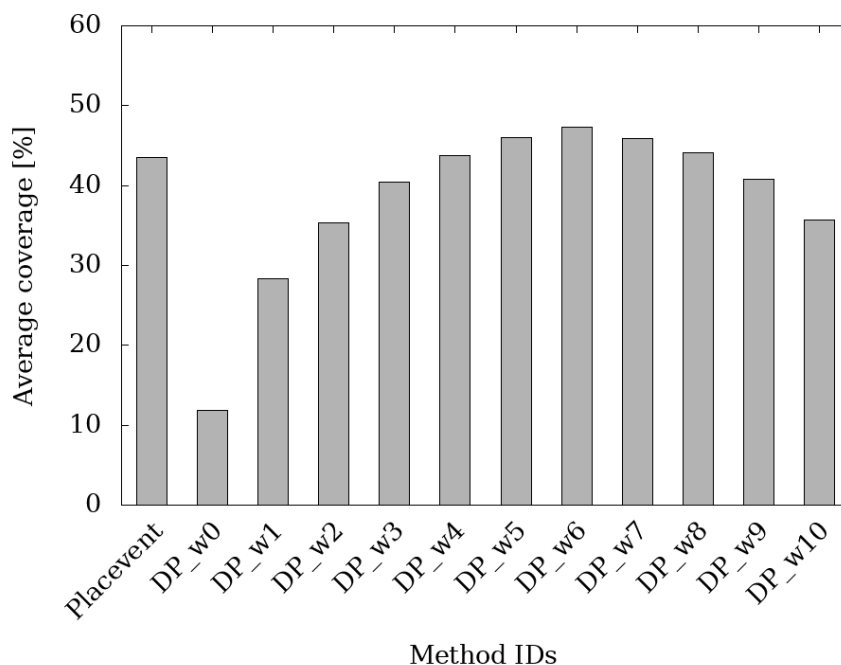


図 19. テストセット全体の結果 (Average coverage)

テストセット全体に対する各手法の予測結果 (Average coverage: 151 ターゲットにおける coverage の中央値の平均値)。横軸は各予測手法 (Method IDs) を、縦軸は Average coverage を示す。DroPred (DP\_w0 から DP\_w10) の Average coverage は、1000 サンプルの coverage の中央値の平均値となる。Placevent からは各ターゲットについてそれぞれ 1 サンプルのみが出力されるため、それらの coverage の平均値が Average coverage となる。

次に、テストセット全体の傾向を確認するため、151 ターゲットの coverage の中央値の平均値 (Average coverage) を算出した (図 19)。DroPred (DP\_w0-DP\_w10) の Average coverage は、151 ターゲットについて、1000 サンプルの coverage の中央値を平均した値である。Placevent からは、1 ターゲットにつき 1 サンプルのみが出力されるため、151 ターゲットについて出力したサンプルの coverage の平均値が Average coverage となる。各手法における coverage の中央値のデータについては、Supporting Information 中の表 S3 に記載した。

個別のターゲットの結果と同様、テストセット全体でも DP\_w1 と比較して DP\_w2 以降の Average coverage が上昇していることから、重みの調整によってサンプルの coverage が上昇することが確認された。DP\_w6 において、最も高い

Average coverage (47.4%) が得られた。Placevent の Average coverage は 43.5% であったことから、DP\_w6 においては、Placevent と同等以上の coverage のサンプルを得ることが期待できると言える。

しかし、DP\_w7 以降では Average coverage が低下する傾向があることから、大きすぎる調整指数の使用は、逆に coverage を悪化させる可能性が高いことも示唆された。

### 3-5. 評価 2 : 分布関数の再現度による評価

#### 3-5-1. 評価の概要

次に、DroPred、即ち重みを調整した重み付きモンテカルロ法によって配置した予測水が、基となった水の三次元分布関数をよく再現しているかを評価した。

DroPred によって配置される予測水が分布関数をよく再現しているのであれば、ある位置に予測水がどの程度配置されるかは、その位置の確率密度の大きさに依存するはずである。つまり、ある位置の確率密度の大きさと、その位置における予測水の出現頻度の間には相関が見られると考えられる。本研究ではこれについて、テストセット中の予測水 (5448 個) が存在する位置で検証を行った。

#### 3-5-2. 結晶水が存在する場所の確率密度

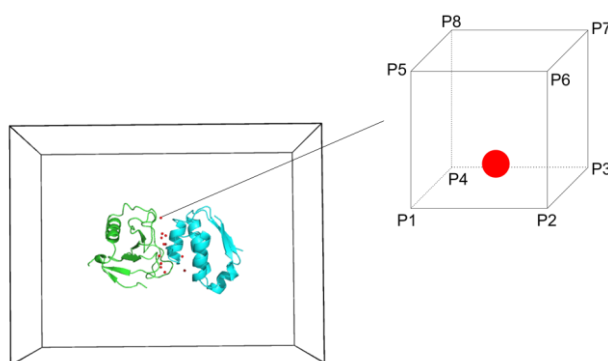


図 20. 結晶水が存在する場所の確率密度

結晶水の座標 (図中赤色の点) は、三次元グリッド空間上の 8 個のグリッド点を頂点としたボクセルのいずれかの内部に存在する。このため本研究では、結晶水が存在する場所の確率密度を、ボクセルの頂点となる 8 つのグリッド点 (P1~P8) における水の確率密度の平均値として定義した。

結晶水の座標は、三次元グリッド空間上の 8 個のグリッド点を頂点とした立方体（ボクセル）の内部に存在することになる（図 20）。このため本研究では、結晶水が存在する場所における確率密度を、結晶水の座標を内包するボクセルの頂点となる 8 つのグリッド点における水の確率密度の平均値として定義した。

### 3-5-3. 予測水の出現率

次に、ある位置に予測水がどの程度出現するかを表す指標として、予測水の出現率（occurrence rate）を定義した。予測水の出現率は、出力したサンプルのうち、その位置から  $1.0\text{\AA}$  以内に予測水を配置することができたサンプルの割合として定義した。これは、検証する位置を結晶水が存在する位置とした場合、その結晶水を再現したサンプルの割合と同義となる。

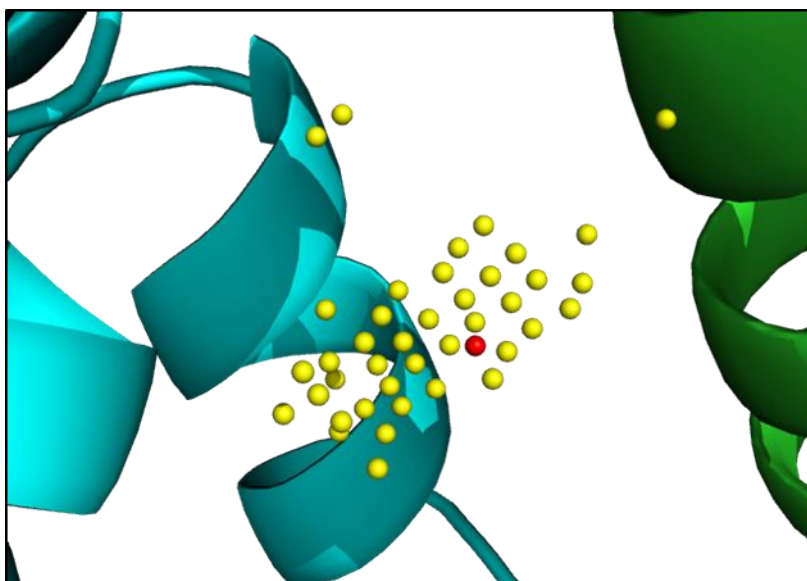


図 21. 結晶水と DroPred による予測水の比較

PDB ID: 3WWT 内のある結晶水を例として、結晶水を赤色の球で、DroPred による予測水を黄色の球で示した。図中の予測水は、1000 個のサンプルのそれぞれにおいてこの結晶水に最も近い位置に配置された予測水を抽出した。予測水の位置は重複する場合がある。結晶水の位置から  $1.0\text{\AA}$  以内に予測水を配置することができたサンプル数を用いて出現率を計算した。

例えば、図 21 に示すターゲット (PDB ID : 3WWT) に含まれるある結晶水は、1000 サンプル中 841 サンプルによって再現されたため、この位置における予測水の出現率は  $841/1000=84.1\%$  と計算される。

#### 3-5-4. 分布関数の再現度についての結果

分布関数の再現度についての結果を示す (図 22)。

一般的なモンテカルロ法を用いた DP\_w0 において、結晶水が存在する場所の確率密度と、予測水の出現率の間の相関係数  $R=0.028$  であり、相関は認められなかった (図 22 (a))。三次元分布関数から得られる確率密度をそのまま重みとして使用した DP\_w1 においては、相関係数  $R=0.647$  と、結晶水が存在する場所の確率密度と、その位置における予測水の出現頻度に相関がみられた (図 22 (b))。重みを調整した DP\_w2 以降では、更に強い相関が確認され、特に、DP\_w3 で相関係数  $R=0.778$  と、最も強い相関がみられた (図 22 (c)~(k))。

この結果から、重み付きモンテカルロ法によって配置された予測水は、基となった水の三次元分布関数をよく表現することが示唆された。

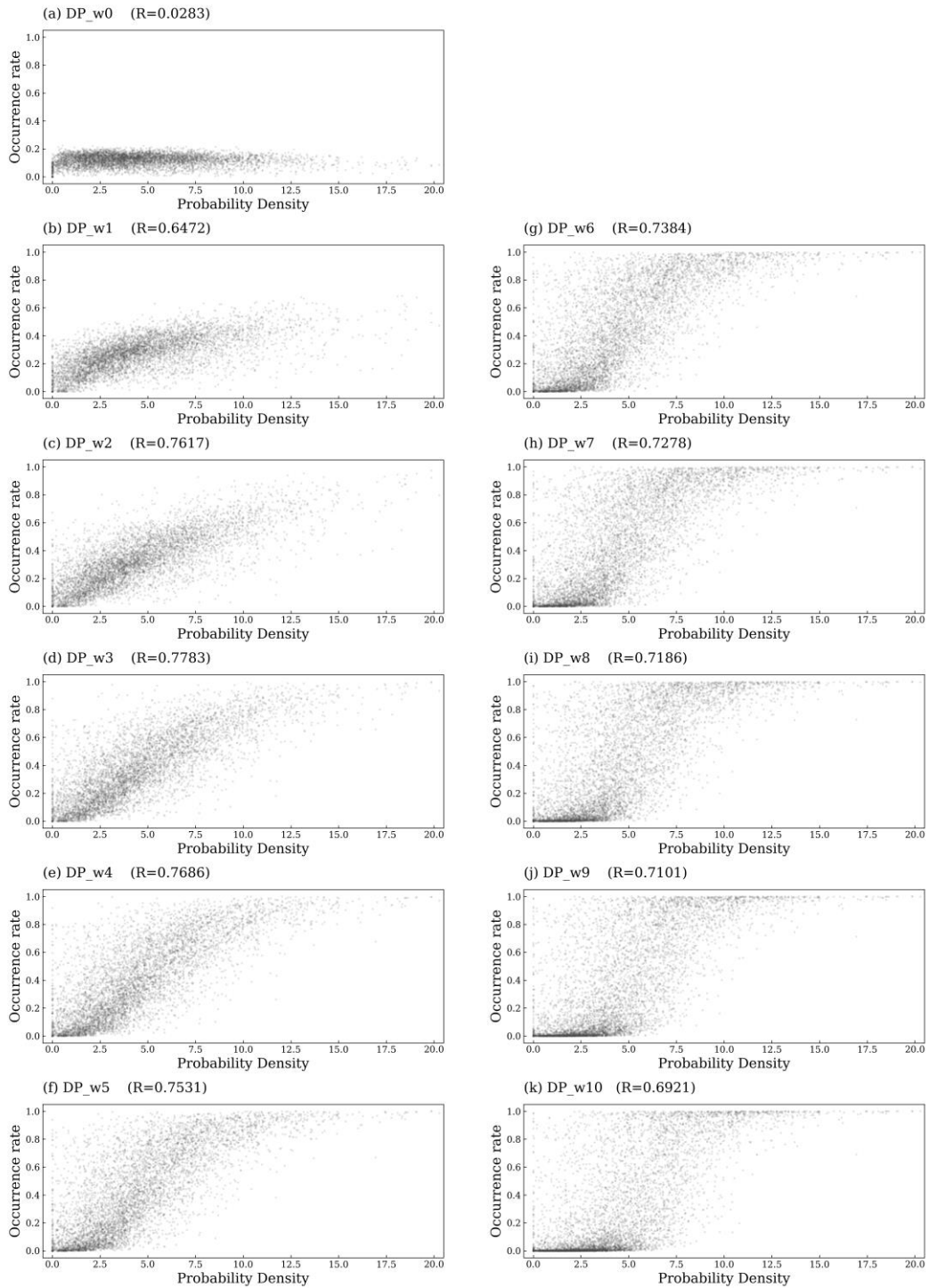


図 22. 結晶水が存在する位置の確率密度と予測水の出現率の相関  
 結晶水が存在する位置の確率密度とその位置における予測水の出現率をプロットした。本研究では、テストセット中で結晶水 (5448 個) が存在する位置について調査を行った。横軸は結晶水が存在する場所の確率密度 ( $g_0$ )、縦軸はその位置における予測水の出現率 (Occurrence rate) をそれぞれ示している。

### 3-6.配置数の制限の影響

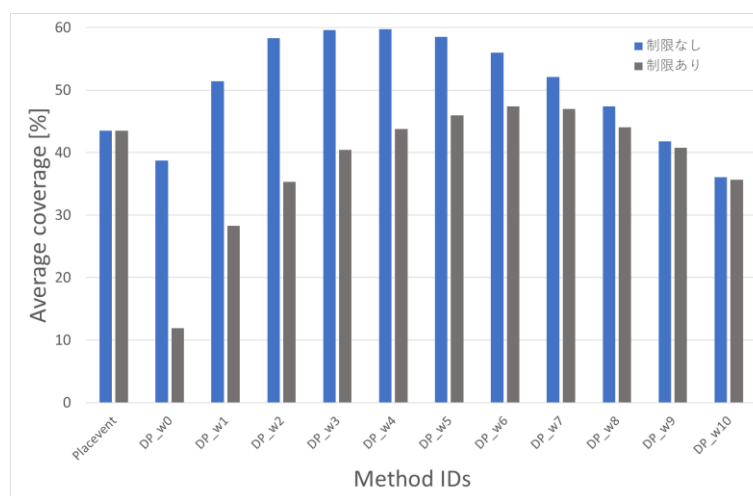


図 23. 配置数の制限の有無による Average coverage の比較

予測水の個数制限の有無による coverage の違いを示す。横軸は予測手法 (Method IDs) を、縦軸は Average coverage を示す。本研究において検証したように、Placevent と同数までの予測水を精度検証に使用した場合の coverage をグレーのバーで、配置した予測水を全て使用した場合の coverage を青色のバーでそれぞれ示した。

本研究では、調整指数毎の DroPred、及び Placevent の結晶水の再現度を比較するために、各サンプルに対して配置した予測水のうち、Placevent と同数までの予測水を評価に使用した。しかし、第 3 章で示したように、DroPred は、使用する調整指数によってはより多くの予測水を配置することが出来る。これについて、配置する水分子の数に制限を設けなかった場合の coverage を計算した(図 23)。この場合、DP\_w4 において最も高い Average coverage (59.8%) が得られた。しかしながら、配置数を制限しない場合 DP\_w4 には平均で 196.8 個の予測水が配置されている。これは Placevent (平均 86.0 個) と比較して多く、実験水を再現していない予測水も多く配置されていると考えられる。

### 3-7. Placevent と DroPred の比較考察

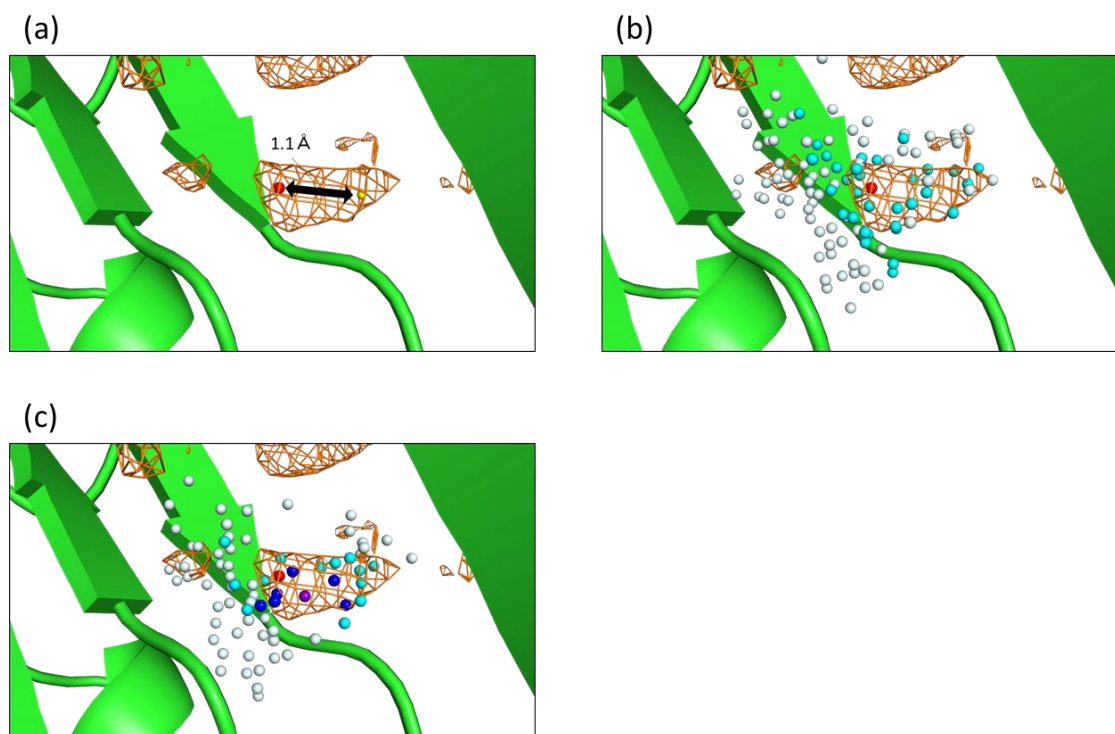


図 24. 各手法によって配置される水分子の比較

タンパク質-タンパク質界面において確率密度が 5.0 以上の領域をオレンジ色のメッシュで示した。また、各手法によって配置した予測水のうち、赤色の球で示した結晶水の最も近くに配置した予測水を図に示している。(a)には Placevent による予測水を黄色の球で示した。(b)は DP\_w1、(c)は DP\_w6 によって配置された予測水を示している。サンプリングにおいては、複数のサンプル間で同じグリッド点に予測水が配置される可能性もあるため、10 未満のサンプルによって配置された予測水を白色の球で、10~49 サンプルによって配置された予測水を水色の球で、50~99 サンプルによって配置された予測水を青色の球で、100 以上のサンプルによって配置された予測水を紫色の球でそれぞれ示している。

ここで、Placevent 及び DroPred による水分子の予測と、結晶水の再現について考察する。図 24 においては、タンパク質-タンパク質界面において確率密度の値が 5.0 以上の領域をオレンジのメッシュで示している。ここに示しているように、同程度の大きさの確率密度の領域にはある程度の広がりがある。

Placevent で配置した予測水を図 24(a)に示す。第 1 章で述べたように、Placevent は確率密度が最も高いグリッド点への配置を繰り返す手法であるため、Placevent が配置する予測水には一定の信頼度があると言える。しかし、最大値と同程度の確率密度を持つ領域に広がりがある場合、結晶水は確率密度が最大のグリッド点とは異なるグリッド点に存在する可能性があり、このような場合には Placevent によって配置した予測水では再現できない可能性がある。実際に、図 25(a)において赤い球で示した結晶水と、黄色の球で示した Placevent の予測水との距離は  $1.1 \text{ \AA}$  であり、この結晶水は Placevent では再現されなかった。

次に、重み付きモンテカルロ法を用いた DP\_w1 によって配置した水分子の分布を示す (図 24(b))。この図では、1000 サンプルそれぞれに配置された予測水のうち、この結晶水に最も近い予測水を示している。サンプリングにおいては、複数のサンプル間で同じグリッド点に水分子が配置される可能性もあるため、同じグリッド点に配置された予測水が 10 サンプル未満の場合は白い球で、10 から 49 サンプルの場合は水色の球で示した。DroPred によるサンプリングでは、複数のグリッド点に予測水を配置することができており、Placevent の予測水と比較して、より結晶水に近いグリッド点にも予測水が存在していることが分かる。このように、重み付きモンテカルロ法を用いたサンプリングを行うことで、確率密度が同程度の領域の広がりに対応し、1つのサンプルでは再現できない結晶水を再現したサンプルを得ることができる可能性があると言える。しかし、多くの予測水がメッシュ外のグリッド点にも配置されていることから分かるように、確率密度が低いグリッド点にも一定数の予測水が配置されてしまうという問題が生じる。DroPred では、既に予測水を配置したグリッド点の周辺のグリッド点には新たに予測水を配置しないため、確率密度の低いグリッド点へ予測水が配置されることで、その周囲に存在する確率密度の高いグリッド点への配置予測水が配置されなくなる。このように、重みの調整を行わない場合、または調整指数が小さな場合には、確率密度の高い場所に存在する結晶水の再現に失敗し、Placevent のような手法と比較してサンプルの coverage が低くなる可能性がある。

更に、調整指数を用いて重みの調整を行った DP\_w6 の予測水の分布を示す (図 24(c))。ここでは先ほどの色分けに加えて、同じグリッド点に配置された予測水が 50 から 99 サンプルの場合は青い球で、100 サンプル以上の場合は紫の球で示している。DP\_w1 に比べて多くの予測水が、確率密度の高いメッシュ内のグリッド点に配置されていることが分かる。これは、重みの調整によって確率密度の高



いグリッド点に予測水がより優先的に配置されるようになったためであると考えられ、重みの調整を行わない場合と比較して確率密度の高いグリッド点への配置が行われなくなる事象が少なくなる可能性が高い。このため、大きな調整指数を使用することで、確率密度の高い場所に存在する予測水をより高確率で再現できるため、サンプルの coverage が高く保たれる可能性が大きくなると言える。

このように、調整指数を用いて、重み付きモンテカルロ法における重みの調整を行うことで、coverage を保ちながら確率密度が同程度の領域の広がりに対応することができると考えられる。

しかしながら、過剰に大きな調整指数を用いた場合には、確率密度の高いグリッド点の優先度が過剰に上昇することで、予測水が確率密度の低いグリッド点へ配置される頻度が低下する可能性がある。実際に、3-3 章で示したように、大きな調整指数を用いた場合には、終了条件を満たすまでに配置される予測水の数が減少するという結果が得られた。これにより、確率密度の低いグリッド点に存在する結晶水の再現度が低下し、coverage が低下する可能性があると考えられる。

## 第4章. 今後の展望

今後の展望として、DroPredによる結晶水の再現度をより高めること、そしてDroPredの適用範囲を広げることに着目して研究を行っていきたいと考えている。具体的には、3D-RISM計算以外の手法による三次元分布関数の使用、DroPredによって得られたサンプルのスコアリング、タンパク質の界面以外の部位や異なる性質を持つ水分子、モデル構造へのDroPredの適用などを考えている。

本研究においては、DroPredによる水分子の位置予測に必要な水の三次元分布関数を3D-RISM計算によって計算したが、DroPredは、「xplor」形式<sup>28</sup>で出力されていれば、どのような手法で計算された三次元分布関数に対しても適用することが可能である。近年では、AIを用いて三次元分布関数を計算する手法なども開発されているため、3D-RISM計算以外の手法で得られた三次元分布関数を基にした予測も行いたいと考えている。

DroPredによる水分子のサンプリングは、1つの三次元分布関数から複数のサンプルを得ることが出来るという特徴がある。結果において、ボックスプロットで示したように、1つのターゲットに対するサンプルであっても、そのcoverageは様々であった。考察でも述べたように複数のサンプルの出力によって、1サンプルでは再現できない結晶水を再現できる可能性が生じ、実際にcoverageがPlaceventを大きく上回るようなサンプルが得られた。しかし一方で、出力した1000サンプルの中には、coverageの低いサンプルも含まれていた。サンプルの出力後、サンプルや配置した水分子1つ1つを評価・予測する手法によって、coverageが高いと考えられるサンプルを選出したり、coverageが低いと考えられるサンプルを棄却したりすることができれば、この手法は更に有用なものとなると考えられる。具体的には、予測水が配置された場所の周辺の物性などを考慮したエネルギー計算などがサンプルの評価のために使用できるのではないかと考えている。

本研究では、タンパク質-タンパク質界面に着目してDroPredの性能評価を行った。今後は、活性部位やリガンド結合部位など、タンパク質-タンパク質界面以外の水分子についても検証を行いたいと考えている。また、本研究ではB-factorが40未満の水分子を評価の対象とした。簡易的な調査では、この基準を満たさない水分子を加えた場合でも、本研究の結果から大きな変化は見られなかったが、これについても今後更なる調査を行いたいと考えている。

コンピュータ計算による水分子の位置予測の利点の1つとして、予測手法によって予測された構造（モデル構造）など、仮想的なタンパク質に対しても適用できるという点が挙げられる。例えば、「AlphaFold」は高精度のタンパク質立体構造予測手法であるが、出力されるモデル構造には水分子が含まれていない。このようなモデル構造は実験によって水分子の位置を決定することが出来ないため、コンピュータ計算による予測手法が有効である。また、モデル構造において、側鎖構造のような細部の精度が十分でない場合には、DroPredのようなサンプリング手法によって水分子の位置を複数提示することが有用である可能性が考えられるため、この点についても検証したいと考えている。

## 第5章. 結論

タンパク質の性質や機能、相互作用を理解する上で、タンパク質周辺の水分子の位置は重要な要素であり、その位置を予測することは有意義であると言える。

本研究では、タンパク質周辺の水分子のサンプリング手法である「DroPred」を開発した。本研究で開発した手法は、三次元分布関数から得られる確率密度を基に予測水を配置するものあり、重み付きモンテカルロ法において使用する重みを調整指数によって調整することで、三次元分布関数をより強く考慮しながら水分子を配置することができる。

本研究では、151 ターゲットからなるタンパク質-タンパク質複合体界面のテストセットに DroPred を適用し、性能評価を行った。

結晶水の再現度による評価では、調整指数を用いて重みを調整することによってサンプリングにおける結晶水の再現度 (coverage) が上昇することが確認された。DroPred から出力したサンプルにおける coverage の中央値を平均した Average coverage は、DP\_w6 において最も高くなり、Placevent によるサンプルの coverage と同程度であった。また、DroPred から得られるサンプルの中には、結晶水の再現度が Placevent を大きく上回るサンプルが含まれていることも確認された。一部のターゲットでは、どの調整指数を用いた場合でも、Average coverage が Placevent を上回らなかったが、そのようなターゲットについても、重みの調整には一定の効果があることが示唆された。

更に、結晶水の存在する位置における確率密度と、その位置における予測水の出現率の相関関係による評価により、重み付きモンテカルロ法によって配置された予測水は、基となった三次元分布関数をよく表現していることも確認された。

本研究の成果は、水分子を含んだタンパク質の機能や相互作用の解析などに貢献できるものであると期待される。

## 謝辞

本研究を遂行するにあたり、北里大学 薬学部 生物分子設計学教室 志鷹真由子教授には指導教官として常に手厚いご指導・ご鞭撻を賜り、また論文の執筆や投稿などにおいても、非常に多くの助言を頂きました。ここに深く感謝の意を表します。

また、北里大学 薬学部 生物分子設計学教室 清田泰臣助教には、プログラムの構築やデータの算出などにおいて、常に有意義なアドバイスを頂きました。深く感謝いたします。

若杉昌輝助教を始めとする北里大学 薬学部 生物分子設計学教室のメンバーには、終始暖かいサポートを頂きましたことを感謝いたします。

最後に、私の博士後期課程までの 9 年間に渡る大学生活を暖かく見守っていただいた両親にも感謝申し上げます。

## 参考文献

1. Berman H.M., Westbrook J., Feng Z., Gilliland G., Bhat T.N., Weissig H., Shindyalov I.N., Bourne P.E., *Nucleic Acids Res*, 28, 235-242 (2000).
2. Samways M.L., Taylor R.D., Bruce Macdonald H.E., Essex J.W., *Chem Soc Rev*, 50, 9104-9120 (2021).
3. Rossato G, Ernst B, Vedani A, Smiesko M., *J Chem Inf Model*, 51, 1867-1881 (2011).
4. Schymkowitz J.W., Rousseau F., Martins I.C., Ferkinghoff-Borg J., Stricher F., Serrano L., *Proc Natl Acad Sci U S A*, 102, 10147-10152 (2005).
5. Sato K., Oide M., Nakasako M., *Sci Rep*, 13, 2183 (2023).
6. Jukič M., Konc J., Gobec S., Janežič D., *J Chem Inf Model*, 57, 3094-3103 (2017).
7. Ross G.A., Morris G.M., Biggin P.C., *PLoS One*, 7, e32036 (2012).
8. Eberhardt J., Forli S., *J Chem Theory Comput*, 19, 2535-2556 (2023).
9. Abel R., Young T., Farid R., Berne B.J., Friesner R.A., *J Am Chem Soc*, 130, 2817-2831 (2008).
10. Jumper J., Evans R., Pritzel A., Green T., Figurnov M., Ronneberger O., Tunyasuvunakool K., Bates R., Židek A., Potapenko A., Bridgland A., Meyer C., Kohl S.A.A., Ballard A.J., Cowie A., Romera-Paredes B., Nikolov S., Jain R., Adler J., Back T., Petersen S., Reiman D., Clancy E., Zielinski M., Steinegger M., Pacholska M., Berghammer T., Bodenstein S., Silver D., Vinyals O., Senior A.W., Kavukcuoglu K., Kohli P., Hassabis D., *Nature*, 596, 583-589 (2021).
11. Evans R., O'Neill M., Pritzel A., Antropova N., Senior A., Green T., Židek A., Bates R., Blackwell S., Yim J., Ronneberger O., Bodenstein S., Zielinski M., Bridgland A., Potapenko A., Cowie A., Tunyasuvunakool K., Jain R., Clancy E., Kohli P., Jumper J., Hassabis D., (BioRxiv) <https://www.biorxiv.org/content/10.1101/2021.10.04.463034v2> (2021).
12. Imai T., Hiraoka R., Kovalenko A., Hirata F., *J Am Chem Soc*, 127, 15334-15335 (2005).
13. Howard J.J., Pettitt B.M., *J Stat Phys*, 145, 441-466 (2011).
14. Roy D., Kovalenko A., *Int J Mol Sci*, 22, 5061 (2021).
15. Imai T., Hiraoka R., Kovalenko A., Hirata F., *Proteins*, 66, 804-813, (2007).

16. Kawama K., Fukushima Y., Ikeguchi M., Ohta M., Yoshidome T., *J Chem Inf Model*, 62, 4460-4473 (2022).
17. Sindhikara D.J., Yoshida N., Hirata F., *J Comput Chem*, 33, 1536-1543 (2012).
18. Fusani L., Wall I., Palmer D., Cortes A., *Bioinformatics*, 34, 1947-1948 (2018).
19. (a) Chiba S., Kiyota Y., Takeda-Shitaka M., Abstract book of the 43rd Symposium on Structural Activity Relationship 2015 & The 10th Japan-China Joint Symposium on Drug Discovery and Development, Niigata, Japan 2015, pp. 74–75. (b) Chiba S., Development of a method for prediction of hydration structures of protein using weighted Monte Carlo method. (2016) [Unpublished master thesis]. Kitasato University.
20. Yakowitz, S., J. E. Krimmel, F. Szidarovszky, *SIAM Journal on Numerical Analysis*, 15, 1289–1300 (1978).
21. Park S., Seok C., *J Chem Inf Model*, 62, 3157-3168 (2022).
22. Krull F., Korff G., Elghobashi-Meinhardt N., Knapp E.W., *J Chem Inf Model*, 55, 1495-1507 (2015).
23. Carugo O., *BMC Bioinformatics*, 19, 61 (2018).
24. Carugo O., *Acta Crystallogr D Struct Biol*, 78, 69-74 (2022).
25. Case D. A., Aktulga H. M., Belfon K., Ben-Shalom I. Y., Berryman J. T., Brozell S. R., Cerutti D. S., Cheatham III T. E., Cisneros G. A., Cruzeiro V. W. D., Darden T. A., Duke R. E., Giambasu G., Gilson M. K., Gohlke H., Goetz A. W., Harris R., Izadi S., Izmailov S. A., Kasavajhala K., Kaymak M. C., King E., Kovalenko A., Kurtzman T., Lee T. S., LeGrand S., Li P., Lin C., Liu J., Luchko T., Luo R., Machado M., Man V., Manathunga M., Merz K. M., Miao Y., Mikhailovskii O., Monard G., Nguyen H., O’Hearn K. A., Onufriev A., Pan F., Pantano S., Qi R., Rahnamoun A., Roe D. R., Roitberg A., Sagui C., Schott-Verdugo S., Shajan A., Shen J., Simmerling C. L., Skrynnikov N. R., Smith J., Swails J., Walker R. C., Wang J., Wang J., Wei H., Wolf R. M., Wu X., Xiong Y., Xue Y., York D. M., Zhao S., Kollman P. A., Amber 2022, University of California, San Francisco (2022).
26. Hornak V., Abel R., Okur A., Strockbine B., Roitberg A., Simmerling C., *Proteins*, 65, 712 (2006).
27. Heo L., Park S., Seok C., *J Chem Inf Model*, 61, 2283-2293 (2021).

28. Brünger A. T., XPLOR manual version 3.0, Yale University (1992).

29. The PyMOL Molecular Graphics System, Version 2.3.0 Schrödinger, LLC.

\*本論文におけるタンパク質の立体構造の図は全て PyMOL で作成した

## 論文目録

本論文は、学術雑誌に掲載された次の論文を基にしたものである。

Kobayashi S., Kiyota Y., Takeda-Shitaka M., Bulletin of the Chemical Society of Japan,97(6), uoae063 (2024). (BCSJ は CSJ と OUP の共同出版)



三次元分布関数を用いた  
タンパク質周辺の水分子のサンプリング手法の開発

補足資料

## 補足資料

### S-1. 検証用テストセット

本研究において、DroPred の検証に使用したテストセットを示す（表 S1）。

表 S1：検証用データセット

本研究において使用したデータセット。左からターゲット番号 (No.)、PDB ID、ターゲットに含まれる chain 名、残基数 (residue)、界面残基数 (residue (int))、ターゲットに含まれる結晶水の個数 (water (all))、B-factor が<sup>s</sup> 40 未満の結晶水の個数 (water (b<40[Å<sup>2</sup>]))、ターゲットの分解能 (resolution[Å])。

No.	PDB ID	chain	residue	residue (IF)	water (all)	water (b<40[Å <sup>2</sup> ])	resolution [Å]
1	1AY7	AB	186	83	28	14	1.7
2	1DX5	MI	377	108	24	24	2.3
3	1DZB	XA	351	115	23	23	2
4	1FYH	AB	443	127	21	9	2.04
5	1GPQ_1	AD	255	110	52	46	1.6
6	1GPQ_2	ABC	382	216	96	90	1.6
7	1HQ3	ABDGH	463	413	100	40	2.15
8	1I2M	AB	553	200	44	30	1.76
9	1IM3	AD	370	94	37	11	2.2
10	1IQD	CAB	564	277	80	69	2
11	1JTD	AB	535	158	49	48	2.3
12	1L4D	AB	361	106	27	19	2.3
13	1LFD	AB	254	79	17	15	2.1
14	1LPB	AB	534	119	22	12	2.46
15	1OAK	AHL	626	284	58	35	2.2
16	1OP9	BA	251	83	27	27	1.86
17	1P2C	ABC	551	292	89	71	2
18	1PXV	CA	294	138	29	21	1.8
19	1TA3	AB	575	155	73	72	1.7
20	1UOS	YA	206	89	27	14	1.9

21	1UUG	AB	304	135	24	15	2.4
22	1V7P	ABC	454	264	114	100	1.9
23	1VFB	ABC	352	169	38	18	1.8
24	1WMH	AB	165	71	31	26	1.5
25	1WRD	AB	174	75	35	26	1.75
26	2ADF	AHL	616	320	82	76	1.9
27	2BKY_1	AXY	277	147	34	29	1.7
28	2BKY_2	YAB	277	178	72	59	1.7
29	2CMR	AHL	604	307	55	46	2
30	2DVW	AB	302	148	48	39	2.3
31	2FHZ	BA	199	138	72	61	1.15
32	2GC7	ABC	611	259	52	42	1.9
33	2GHW	AB	423	147	26	26	2.3
34	2HQS	AH	528	168	64	54	1.5
35	2I25	LN	243	97	25	11	1.8
36	2IBG	AF	346	87	11	8	2.2
37	2ID0	AB	248	123	39	26	2.1
38	2JBG	BAC	293	155	56	48	2.2
39	2NQD	AB	332	130	71	63	1.75
40	2OZN	AB	264	94	24	20	1.6
41	2PTT	AB	211	97	35	32	1.63
42	2UYZ	AB	270	87	50	44	1.4
43	2VDU	FD	577	122	19	19	2.4
44	2VXQ	AHL	518	291	96	61	1.9
45	2WBW	AB	315	117	57	46	1.55
46	2WY3	AB	325	135	44	26	1.8
47	2WY7	AQ	367	91	35	28	1.7
48	2WY8	AQ	372	111	47	29	1.7
49	2XGY	AB	294	102	41	15	1.8
50	2XWT	ABC	662	351	131	116	1.9
51	2Y1L	CE	305	152	43	35	1.8
52	2YC1	CAB	297	158	64	54	1.9
53	3CHW	AP	512	139	46	45	2.3
54	3CQX	BC	458	105	23	22	2.3

55	3F1P	AB	251	125	43	43	1.17
56	3F62	BA	253	112	20	14	2
57	3GCG	AB	325	163	47	37	2.3
58	3KF6	BA	241	111	37	18	1.65
59	3KLD	AB	638	118	28	18	2
60	3L9J	TC	279	89	29	19	2.1
61	3M18	AB	372	155	52	43	1.95
62	3MA9	AHL	623	323	82	54	2.05
63	3O2D	AHL	616	330	62	62	2.19
64	3P9W	AB	219	101	31	14	2.41
65	3REP	BA	392	105	29	20	1.8
66	3RJ3	AD	422	126	20	10	2.35
67	3RKD	AHL	577	310	87	68	1.9
68	3RNK	AB	218	119	15	12	1.74
69	3TDZ	AC	269	94	40	40	2
70	3U30	ABC	578	311	36	33	2.43
71	3UFX	AB	666	268	56	40	2.35
72	3V60	CEA	613	306	82	45	1.95
73	3W9E	ABC	655	326	51	7	2.3
74	3WDG	AB	326	134	44	44	2.2
75	3WIH	AHL	526	307	108	65	1.7
76	3WWT	AB	231	113	33	21	2
77	3ZDM	ABC	198	126	50	29	1.8
78	4A5U	AB	240	94	24	13	2
79	4AG1	AC	294	126	39	32	1.4
80	4BI8	AB	259	123	37	35	2
81	4BL7	BA	375	163	43	41	1.89
82	4CMH	ABC	686	338	117	71	1.53
83	4CZX	AB	438	99	30	17	1.85
84	4DCK	AC	301	139	37	8	2.2
85	4DTG	KHL	520	312	98	39	1.8
86	4E5X	ABG	474	253	67	45	1.95
87	4G6M	AHL	583	307	93	57	1.81
88	4G7X	AB	216	105	35	35	1.44

89	4H8W	GHL	777	296	87	33	1.85
90	4HCR	AHL	635	309	106	88	2.3
91	4HEM	BCE	387	239	112	57	1.65
92	4HPL	AB	196	130	23	22	2
93	4IMI	AB	524	170	35	10	2.35
94	4IOI	BH	284	83	34	17	1.95
95	4J7B	AB	497	143	44	30	2.3
96	4JE4	AB	195	110	21	21	2.31
97	4JHP	BC	506	142	46	31	1.9
98	4K12	BA	163	79	49	48	1.08
99	4KT6	BA	415	189	72	63	1.71
100	4L5N	ACD	329	153	50	24	2.16
101	4LGR	AB	379	92	25	17	1.65
102	4M1G	BHL	510	285	85	85	1.6
103	4MA7	AHL	542	289	79	30	1.97
104	4N6R	AB	354	184	47	43	2.2
105	4N90	AB	233	83	25	25	1.5
106	4NBX	AB	266	97	33	28	1.75
107	4NBY	AC	371	94	37	23	2.08
108	4NRH	AD	459	125	24	7	2.2
109	4P78	AD	152	119	16	13	2.12
110	4PJ2	AD	272	105	37	27	1.24
111	4QAF	CDA	305	207	64	46	1.8
112	4QLP	AB	390	227	122	119	1.1
113	4RGO	SHL	662	317	122	51	1.8
114	4TSB	AHL	558	283	71	67	1.95
115	4TXV	AB	317	131	37	25	2
116	4UHP	CD	215	99	22	18	2
117	4WEM	AB	375	117	48	34	1.55
118	4WKZ	BA	424	129	29	24	1.79
119	4YDY	AI	280	140	28	20	2
120	5A6W	BC	159	99	29	28	1.6
121	5BV7	AHL	809	320	83	27	2.45
122	5BVP	HLI	581	324	76	30	2.2

123	5C50	AB	388	175	41	29	1.63
124	5CEC	AB	598	190	66	62	1.36
125	5CZX	AHL	655	305	93	57	2.1
126	5D2M	AC	319	86	27	24	2.4
127	5D93	ABC	659	286	67	62	2.2
128	5DC4	AB	192	110	31	24	1.48
129	5DJT	AB	209	103	37	32	1.4
130	5D0I	AE	246	102	14	13	2.2
131	5EE4	ACD	584	230	29	25	2.3
132	5ELU	AB	166	81	13	12	2.35
133	5F1K	AC	368	111	23	21	2.3
134	5F72	KS	511	157	25	25	1.85
135	5GGS	ABZ	549	325	86	57	2
136	5H5Z	BA	376	204	73	34	1.74
137	5IWB	AB	229	125	41	32	1.76
138	5J3T	ACB	395	218	48	17	1.6
139	5JDS	AB	245	84	20	19	1.7
140	5KVF	EHL	551	314	132	98	1.4
141	5KW9	AHL	723	300	53	20	2.3
142	5L21	AB	544	162	84	41	1.68
143	5LB7	BA	248	144	47	24	1.5
144	5M20	AB	218	94	33	30	1.26
145	5002	AC	452	74	24	18	1.72
146	5V1Y	AC	191	67	36	25	1.42
147	5VAG	BCA	665	328	82	18	1.9
148	5WCA	HL	442	242	123	111	1.37
149	5WV0	BC	311	101	30	26	2
150	6APP	AB	182	75	39	36	1.75
151	6AZZ	DEF	599	308	37	21	2.4

## S-2. 各手法における水分子の配置数

各予測手法において、配置が終了するまでに配置された予測水の数を示す(表 S2)。

表 S2: 水分子の配置数

各ターゲットに対する予測水の配置数を示す。DroPred (DP\_w0 から DP\_w10) の配置数は 1000 サンプルの平均値とした。

PDB ID	Placevent	DP_w0	DP_w1	DP_w2	DP_w3	DP_w4	DP_w5	DP_w6	DP_w7	DP_w8	DP_w9	DP_w10
1AY7	45	120.7	92.2	107.1	99.4	89.2	76.8	65.2	54.0	44.3	36.5	30.3
1DX5	63	192.7	146.0	172.9	162.0	145.1	124.5	103.5	82.8	64.8	50.1	38.8
1DZB	57	168.6	128.4	154.7	142.9	125.7	105.3	83.8	64.4	51.2	41.6	34.9
1FYH	58	218.0	163.7	193.0	180.8	161.8	137.8	111.8	90.5	71.8	53.7	39.9
1GPQ_1	69	184.5	135.9	159.7	148.1	129.6	107.1	83.6	60.5	42.3	30.3	22.4
1GPQ_2	124	360.0	270.6	320.3	295.9	262.2	222.2	179.5	138.4	103.3	79.1	60.9
1HQ3	71	964.4	737.3	891.4	822.5	725.8	604.4	477.7	361.4	261.6	180.5	123.0
1I2M	76	267.6	210.4	244.9	230.1	210.4	184.6	158.3	133.0	112.2	94.8	82.0
1IM3	64	167.8	130.9	153.9	144.7	129.7	111.3	92.7	75.0	62.0	51.2	42.4
1IQD	98	341.3	250.2	297.4	276.7	246.9	212.6	175.9	144.9	117.9	95.3	76.2
1JTD	52	223.9	171.7	201.1	184.3	161.4	135.4	109.7	85.6	67.0	53.2	42.9
1L4D	54	148.4	109.6	128.9	118.5	105.1	89.2	73.9	58.7	44.4	33.8	27.8
1LFD	43	121.1	90.4	105.8	98.1	86.2	71.2	56.7	43.9	32.3	23.4	17.6
1LPB	76	203.2	151.9	174.7	163.9	145.8	116.9	84.4	56.5	36.1	22.2	11.7
1OAK	111	319.7	240.9	285.5	263.8	232.2	193.5	154.2	122.5	97.3	75.5	59.1
1OP9	37	110.6	79.9	93.7	84.9	73.3	60.0	48.1	38.4	30.8	24.6	20.2
1P2C	135	372.1	282.3	333.4	303.8	261.2	205.9	149.0	101.3	66.8	42.7	27.4
1PXV	74	226.6	170.9	200.9	183.5	157.2	121.6	83.0	50.1	29.5	17.6	10.8
1TA3	91	228.5	178.9	211.8	198.0	178.4	157.6	135.5	115.1	95.4	79.6	65.9
1U0S	56	140.4	109.1	128.4	120.1	108.2	92.4	76.0	60.2	46.7	35.5	26.7
1UUG	48	222.7	168.8	199.7	186.7	166.4	140.4	111.1	84.8	64.2	46.4	32.1
1V7P	155	408.2	309.2	366.6	340.3	301.4	256.8	213.4	172.5	137.7	109.2	87.7
1VFB	123	298.4	227.9	266.9	250.3	226.2	195.1	159.0	123.0	90.0	64.2	46.7
1WMH	43	118.5	86.4	102.8	93.6	79.9	64.5	49.7	36.4	25.7	18.4	13.6

1WRD	51	144.5	106.2	124.7	117.2	106.8	94.3	82.0	70.2	59.7	50.7	43.7
2ADF	155	432.3	333.7	397.6	367.2	322.6	271.2	216.6	167.8	127.0	96.0	72.8
2BKY_1	86	235.7	175.7	209.0	194.7	169.6	135.4	99.7	71.2	48.7	30.4	18.4
2BKY_2	123	332.4	247.4	289.0	267.5	233.1	188.7	144.5	106.9	76.2	51.3	34.2
2CMR	69	449.7	342.5	405.6	375.9	327.7	264.2	197.4	138.5	86.6	49.3	27.7
2DVW	83	214.2	161.1	189.7	172.6	148.5	119.8	94.1	73.9	55.6	39.4	26.2
2FHZ	85	210.7	166.9	193.3	176.3	151.4	123.3	98.2	77.4	59.6	44.2	32.4
2GC7	49	357.1	275.9	326.6	302.6	267.2	224.9	183.2	143.4	111.4	86.8	68.5
2GHW	66	205.9	156.7	180.9	165.9	144.3	120.2	96.1	77.6	64.5	54.5	47.8
2HQS	42	196.1	156.7	184.4	170.4	146.4	115.3	79.7	52.0	34.6	22.0	13.1
2I25	47	126.6	93.1	108.6	100.7	89.4	74.9	58.2	40.7	27.4	19.0	12.1
2IBG	53	156.9	117.0	136.5	125.8	108.6	85.5	65.5	48.2	33.4	24.5	18.6
2IDO	4	176.0	131.2	153.4	143.5	128.8	110.0	89.1	69.3	52.7	38.5	26.8
2JBG	112	310.1	239.1	282.7	258.5	223.1	175.3	125.1	83.4	55.1	35.5	22.7
2NQD	64	178.2	134.1	157.6	147.4	132.7	114.1	95.7	79.7	66.6	54.6	45.1
2OZN	44	132.0	97.9	113.5	105.3	93.8	80.2	66.9	54.5	43.4	33.6	26.1
2PTT	62	168.6	130.6	151.4	138.9	120.8	99.4	76.7	57.5	43.0	32.1	24.4
2UYZ	62	164.9	121.6	142.5	133.9	121.3	105.0	87.3	70.9	55.1	41.1	31.1
2VDU	69	192.7	144.2	170.1	159.5	144.7	127.8	109.4	92.2	75.9	61.6	49.5
2VXQ	147	423.5	319.6	378.9	350.8	310.8	264.8	219.1	177.4	141.4	109.9	86.2
2WBW	75	211.8	155.6	182.9	168.8	148.1	125.5	101.3	78.2	60.7	46.2	34.8
2WY3	85	228.8	176.1	206.3	192.5	173.4	148.6	124.8	102.0	78.8	59.6	45.7
2WY7	50	158.3	109.8	129.3	120.0	107.0	90.5	74.8	61.3	50.4	41.7	34.6
2WY8	78	191.5	150.1	175.1	164.0	148.2	127.2	103.7	81.1	63.5	50.7	42.2
2XGY	57	150.4	114.3	132.5	123.0	108.4	91.2	74.4	59.6	47.8	37.2	29.0
2XWT	72	467.0	351.4	417.5	385.4	337.6	281.6	228.6	181.2	142.9	112.4	90.7
2Y1L	58	169.9	132.5	153.8	141.1	122.0	97.4	73.4	54.2	38.2	26.5	16.8
2YC1	109	284.5	219.9	260.8	243.5	215.3	178.1	138.1	104.3	75.2	51.3	37.1
3CHW	44	194.3	148.7	175.1	163.9	147.3	128.5	110.4	92.3	74.1	59.6	48.6
3CQX	71	193.9	143.1	170.2	158.3	141.9	121.8	100.4	79.8	61.9	47.9	38.1
3F1P	90	262.4	191.8	221.7	207.4	183.4	151.8	120.1	96.5	77.3	62.1	49.9
3F62	57	168.0	135.4	158.3	147.0	130.3	108.3	83.6	62.8	46.5	35.3	27.0
3GCG	100	251.7	191.7	222.5	206.0	179.9	151.0	122.9	101.5	83.0	67.0	53.6
3KF6	59	184.1	140.4	164.2	154.8	142.1	125.2	107.6	89.8	72.8	58.1	46.7



3KLD	53	161.8	116.9	135.3	123.3	105.5	84.0	59.7	42.5	31.3	22.8	16.2
3L9J	63	168.3	123.8	145.3	134.9	120.1	101.9	83.3	67.5	55.4	45.4	37.3
3M18	82	206.9	162.9	190.5	176.7	154.9	126.6	96.1	68.1	48.6	34.5	25.0
3MA9	59	414.9	310.2	366.9	339.1	296.6	246.5	196.4	153.2	117.8	90.7	69.7
3O2D	193	503.8	382.1	456.1	420.4	371.5	314.4	260.0	206.1	156.0	115.2	84.2
3P9W	62	179.6	136.9	160.2	148.9	131.0	110.7	91.5	73.4	57.1	43.0	32.3
3REP	62	162.6	125.0	147.7	139.0	125.3	108.0	90.0	70.6	52.8	39.7	31.2
3RJ3	71	193.2	149.5	173.5	158.8	137.5	110.3	83.9	62.0	45.9	34.2	24.9
3RKD	159	451.4	344.9	414.0	381.4	336.3	280.7	222.4	166.0	120.0	82.5	57.5
3RNK	62	170.0	128.9	147.4	133.9	114.3	93.6	77.0	61.5	47.3	35.9	27.1
3TDZ	61	163.0	123.9	143.6	135.1	121.4	104.4	86.6	70.5	57.0	44.4	34.6
3U30	140	455.2	340.6	404.6	374.9	334.1	283.3	234.9	194.6	158.0	123.3	94.5
3UFX	144	410.6	323.8	384.4	356.8	318.4	272.5	226.0	181.9	144.2	113.9	90.2
3V60	167	456.0	343.1	409.3	379.7	334.9	283.2	228.7	179.0	136.6	100.2	70.3
3W9E	167	501.4	379.5	452.8	415.6	362.2	296.3	231.6	171.2	119.8	78.7	53.0
3WDG	84	221.2	170.3	198.4	184.7	164.9	139.8	114.3	92.5	74.4	60.8	50.5
3WIH	176	470.7	354.5	420.3	391.5	352.7	305.1	252.5	202.6	159.5	126.7	98.8
3WWT	64	170.5	131.7	152.8	142.6	127.9	109.8	89.6	72.0	56.8	45.4	36.4
3ZDM	85	224.7	174.4	208.2	195.3	176.2	152.0	126.7	100.7	77.3	60.6	48.7
4A5U	61	156.1	116.6	135.6	125.0	109.7	91.9	75.1	59.2	45.5	35.7	29.1
4AG1	85	209.8	161.7	187.1	173.3	151.4	124.3	96.6	72.4	51.8	36.0	25.4
4B18	67	186.3	139.4	166.1	154.1	138.6	120.8	103.0	84.7	67.3	51.7	37.6
4BL7	96	267.7	199.7	233.9	218.0	195.3	167.6	139.3	115.1	94.2	76.8	63.0
4CMH	163	464.7	349.5	417.6	386.9	344.4	294.5	246.2	203.1	165.4	133.2	105.4
4CZX	33	135.5	100.9	116.5	107.0	94.9	81.8	68.5	55.5	45.7	37.9	31.3
4DCK	75	193.7	147.4	174.3	161.5	143.2	121.0	97.8	76.4	59.6	46.5	36.7
4DTG	184	505.7	382.4	459.0	421.1	366.8	301.8	238.3	183.8	136.2	95.4	66.1
4E5X	154	424.3	328.9	391.9	361.5	314.2	255.9	197.7	147.2	108.1	81.1	61.3
4G6M	131	415.9	309.7	371.7	343.3	306.7	263.5	221.7	184.7	153.7	129.3	108.6
4G7X	68	197.8	148.1	172.6	160.4	142.0	120.6	98.9	80.9	66.7	56.5	48.8
4H8W	42	333.6	247.8	292.5	269.7	233.4	189.1	144.3	106.2	74.9	52.0	37.0
4HCR	154	406.0	296.3	348.0	324.6	288.7	243.7	193.4	149.0	115.9	90.0	68.3
4HEM	204	570.9	424.3	503.3	461.7	399.3	316.4	231.5	165.3	118.2	85.4	63.2
4HPL	85	244.4	190.9	226.0	208.6	185.4	156.5	126.8	97.1	70.7	51.4	38.9

4IMI	80	256.5	200.7	235.0	217.9	194.2	167.6	141.7	117.8	96.0	76.7	61.2
4I01	38	109.8	82.7	95.3	87.8	76.9	65.4	54.1	43.4	34.6	27.7	22.7
4J7B	86	236.9	182.5	217.3	202.6	183.2	161.9	140.1	119.7	99.9	82.8	68.5
4JE4	63	170.1	128.7	152.2	141.7	126.5	109.5	92.4	78.2	67.5	58.5	50.7
4JHP	71	196.4	150.0	177.7	166.2	150.7	132.3	114.4	97.5	81.5	65.6	51.8
4K12	36	150.0	114.1	134.9	127.7	117.2	102.8	85.8	70.0	56.5	43.9	32.2
4KT6	97	298.8	231.4	270.3	252.6	224.0	189.0	153.2	114.3	77.6	50.6	33.9
4L5N	69	278.7	213.6	252.7	236.6	211.8	182.1	154.9	132.2	112.7	94.5	78.9
4LGR	38	122.0	91.9	106.5	98.1	85.6	71.0	55.4	42.1	33.7	27.5	23.2
4M1G	128	429.8	317.8	380.3	351.1	312.4	268.5	222.8	180.9	145.7	114.1	87.8
4MA7	137	370.8	279.5	335.6	311.3	277.2	238.0	200.4	166.1	135.6	112.2	92.1
4N6R	57	298.1	230.7	271.7	252.8	227.4	193.9	159.0	127.2	100.5	78.6	62.8
4N90	48	130.9	102.3	118.5	110.3	99.0	85.7	72.3	61.0	49.6	40.5	33.7
4NBX	53	157.2	118.5	139.3	130.1	116.4	100.2	83.3	66.7	52.6	42.9	35.9
4NBY	67	173.3	128.8	151.2	140.9	125.0	105.0	83.9	65.0	49.2	36.3	27.0
4NRH	68	178.0	132.3	157.9	145.5	128.6	108.2	88.4	69.6	55.8	45.8	37.1
4P78	29	202.1	152.2	178.7	165.3	145.8	122.4	99.8	76.7	57.0	42.6	31.8
4PJ2	68	188.4	137.4	158.5	148.4	134.1	115.5	93.8	72.1	53.9	39.8	28.7
4QAF	125	373.8	279.4	330.3	306.9	272.3	227.3	181.0	135.3	98.8	72.9	52.5
4QLP	83	386.4	295.4	344.7	321.2	284.9	236.4	184.3	138.3	104.3	77.8	59.6
4RGO	54	407.1	307.1	368.0	341.9	302.8	257.4	210.6	166.2	125.9	92.3	70.0
4TSB	142	428.0	318.7	378.2	347.8	305.8	251.2	192.9	141.4	100.9	73.2	54.5
4TXV	65	195.2	144.7	168.8	156.7	137.4	114.4	89.3	67.1	53.0	42.2	33.7
4UHP	24	139.2	107.3	125.8	116.9	103.7	86.7	68.9	54.4	45.5	39.8	34.7
4WEM	62	182.1	137.9	161.4	148.4	128.5	105.1	82.4	64.7	52.4	43.5	36.7
4WKZ	59	160.5	117.3	137.3	127.5	112.3	93.9	75.3	59.1	46.2	36.4	29.4
4YDY	71	193.1	150.6	175.1	162.2	143.0	121.2	100.1	79.6	62.7	51.0	42.3
5A6W	65	174.2	127.8	148.4	138.3	124.0	107.9	92.5	78.0	65.4	53.7	43.5
5BV7	127	446.5	332.2	396.8	366.3	318.1	257.7	193.3	137.3	91.8	61.9	39.8
5BVP	180	491.2	384.4	461.2	428.3	379.3	322.5	263.5	209.4	163.5	128.4	100.4
5C50	89	323.6	238.6	278.4	263.2	239.4	210.8	180.8	153.3	130.8	112.3	97.6
5CEC	94	285.1	217.3	255.6	236.4	207.0	169.4	129.5	95.4	71.0	53.6	40.2
5CZX	159	484.5	363.6	435.9	405.6	359.6	303.3	244.0	188.8	145.0	108.6	81.2
5D2M	62	172.3	126.1	146.4	138.0	126.3	110.1	91.7	74.0	59.5	46.9	36.4

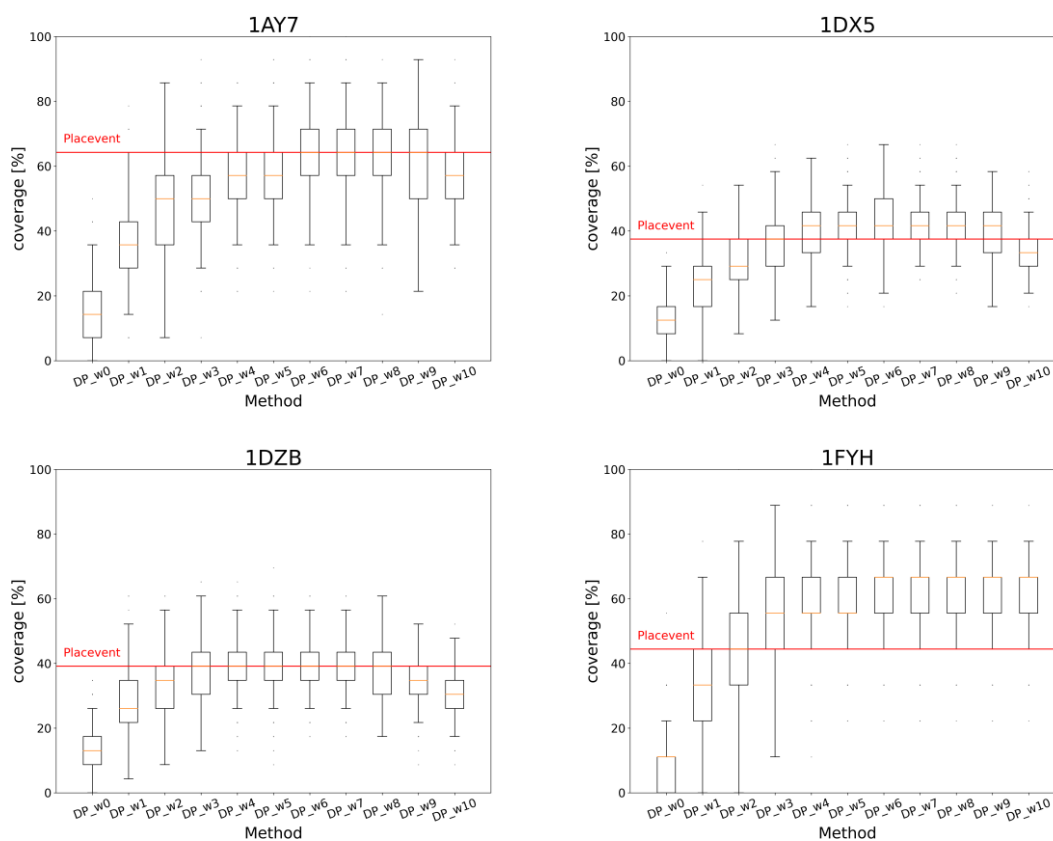
5D93	142	410.1	305.2	362.0	336.7	301.1	260.6	218.2	183.0	156.4	133.8	113.5
5DC4	55	155.2	117.2	137.5	127.0	110.9	89.5	67.2	49.0	35.6	25.1	17.5
5DJT	40	137.3	102.8	122.8	113.1	100.1	84.4	68.5	54.2	42.5	33.5	26.5
5D0I	42	151.7	112.3	131.3	122.8	109.5	92.0	74.8	60.4	46.0	34.6	25.4
5EE4	145	431.6	333.8	398.6	370.2	331.6	284.7	233.4	182.9	139.2	104.4	75.7
5ELU	47	141.6	110.1	127.2	118.4	104.5	86.6	69.6	55.1	43.3	33.6	26.1
5F1K	56	166.2	130.0	151.8	139.8	122.9	103.5	85.8	70.2	56.7	46.0	38.2
5F72	87	219.1	165.8	197.7	184.0	165.0	143.1	119.0	93.6	69.3	49.1	32.9
5GGS	169	461.7	351.6	425.9	392.2	346.7	292.6	240.8	198.3	163.6	133.3	107.1
5H5Z	121	323.2	249.6	296.0	272.8	241.2	203.2	163.9	128.6	97.2	74.3	58.8
5IWB	61	166.6	130.6	153.2	141.7	125.5	106.1	88.4	73.7	61.1	50.4	41.9
5J3T	82	346.1	256.9	304.3	281.7	249.9	209.0	164.6	123.7	92.8	68.9	48.8
5JDS	44	128.9	95.0	110.4	102.5	91.5	77.8	65.1	54.5	44.1	35.4	29.0
5KVF	151	468.8	355.0	421.8	387.1	337.5	278.1	219.0	164.9	117.7	82.2	59.3
5KW9	129	351.0	261.2	309.6	285.1	252.1	212.9	174.0	140.8	114.5	93.9	76.3
5L2I	98	270.9	214.2	253.4	239.3	218.8	193.9	169.8	147.3	126.8	107.9	92.2
5LB7	69	220.9	168.1	197.7	185.9	168.8	145.6	117.6	87.8	58.6	34.5	19.7
5M2O	41	138.4	101.8	119.1	108.5	92.3	74.0	55.2	37.6	26.5	20.3	16.8
5002	38	107.0	80.0	93.9	88.2	79.0	66.8	55.3	45.7	37.2	31.2	26.2
5V1Y	69	166.6	127.1	147.6	140.0	127.4	110.0	89.6	72.5	58.8	47.1	37.1
5VAG	160	517.0	400.7	480.1	444.4	392.8	329.2	268.6	219.1	178.1	142.1	112.3
5WCA	154	454.2	312.5	369.3	346.0	311.2	268.0	218.3	172.1	137.3	110.2	88.6
5WVO	64	181.5	140.2	166.2	155.8	141.6	124.0	105.8	89.6	75.3	61.0	48.1
6APP	47	123.4	95.5	111.1	103.2	91.8	77.9	65.9	55.5	48.2	42.4	37.9
6AZZ	152	461.7	356.2	420.9	388.3	340.8	285.8	230.2	178.7	141.2	114.4	94.8

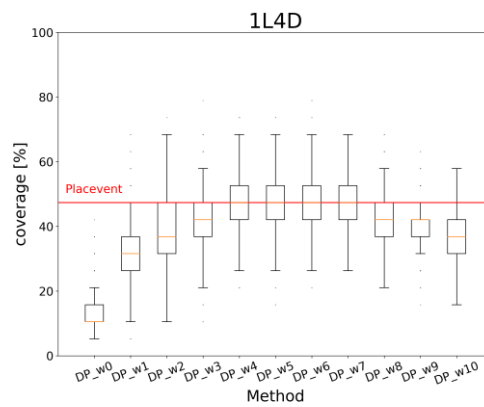
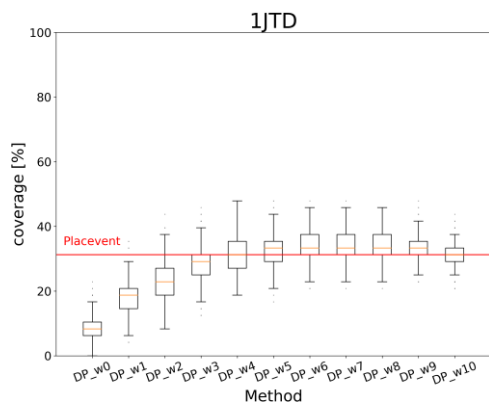
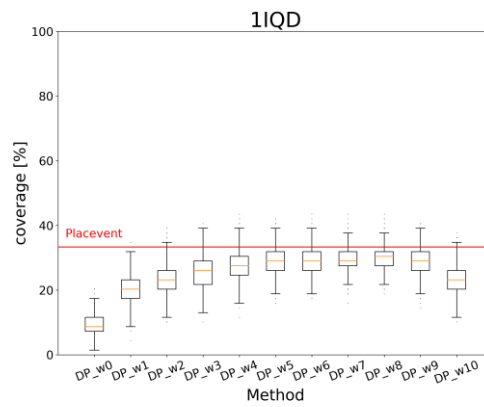
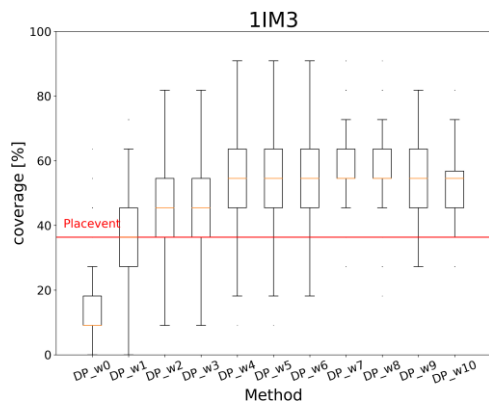
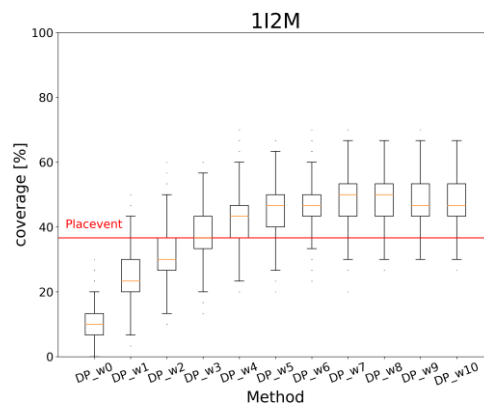
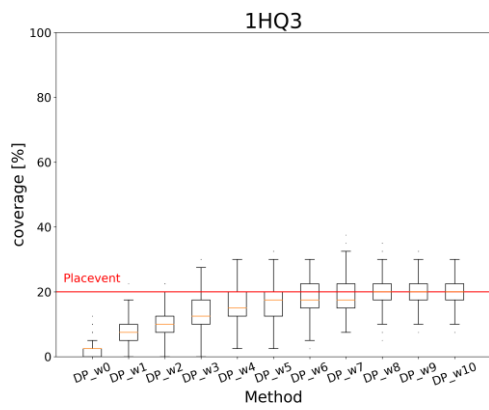
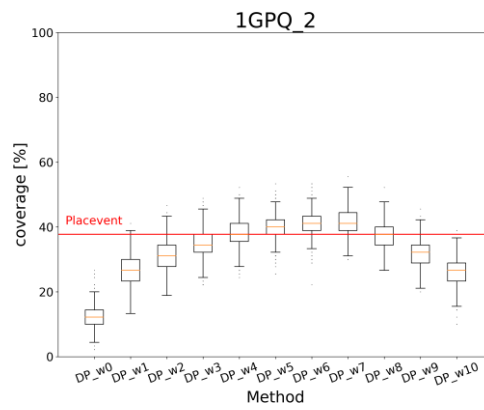
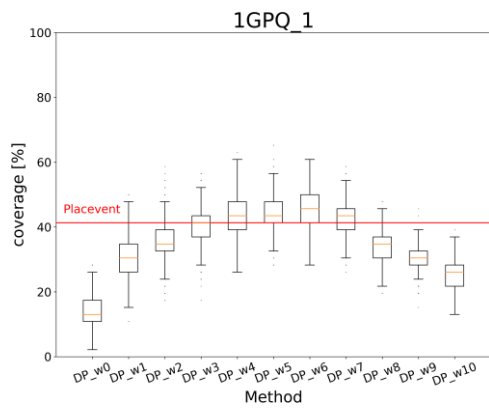
### S-3. ターゲット毎の検証結果

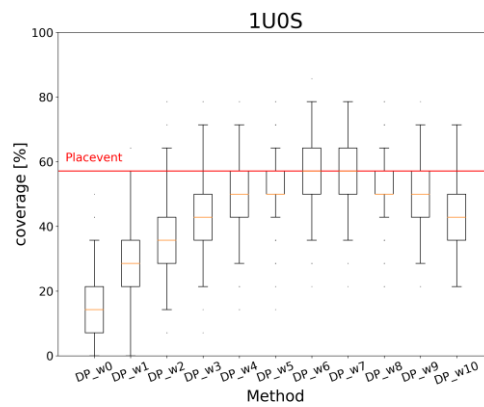
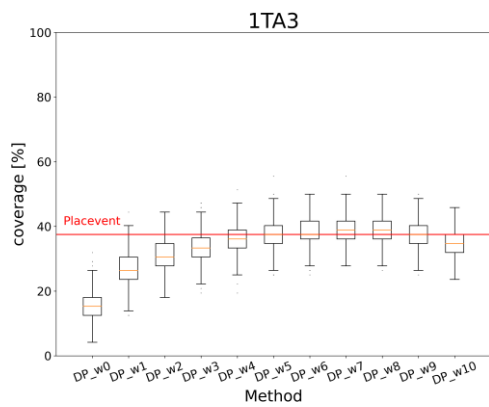
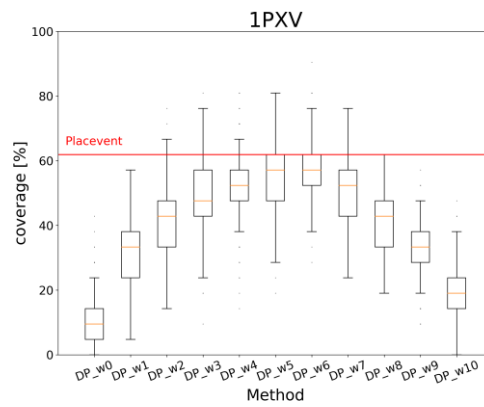
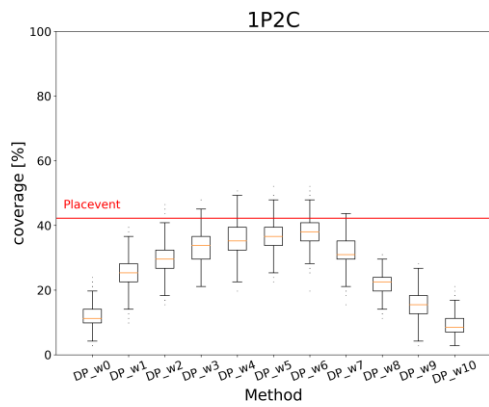
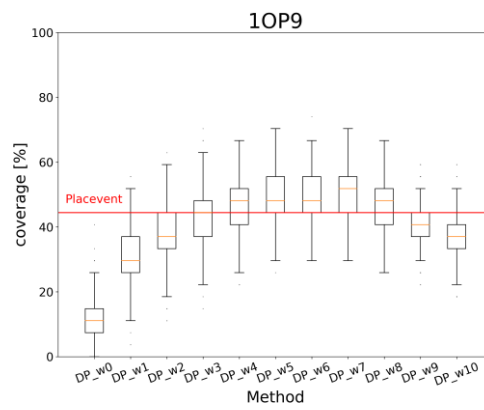
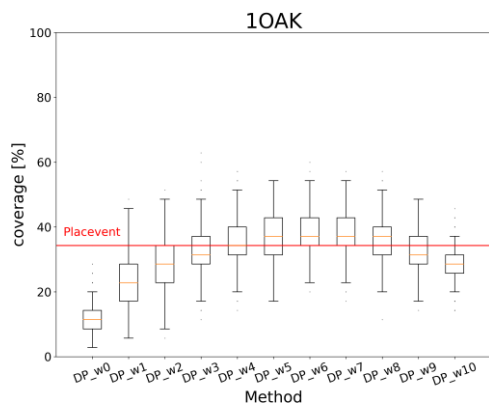
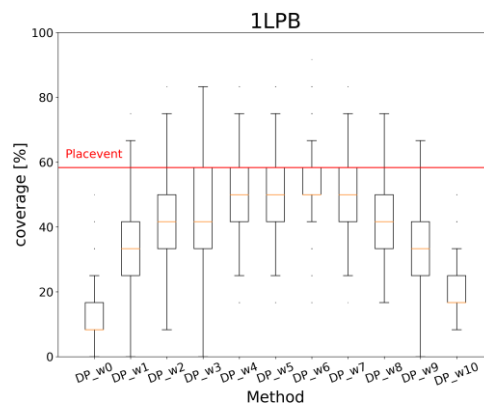
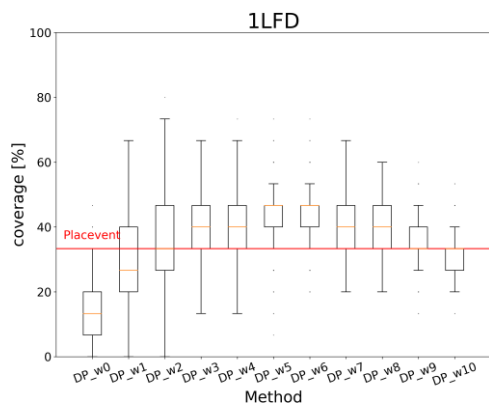
テストセットに含まれる 151 ターゲットについて結晶水の再現度について評価を行った結果を示す (図 S1)。

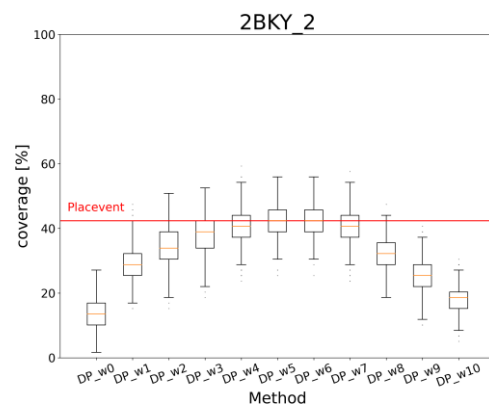
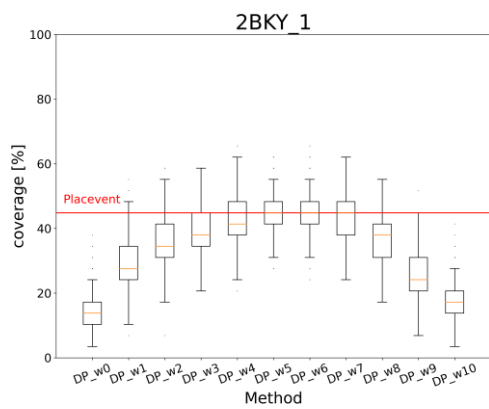
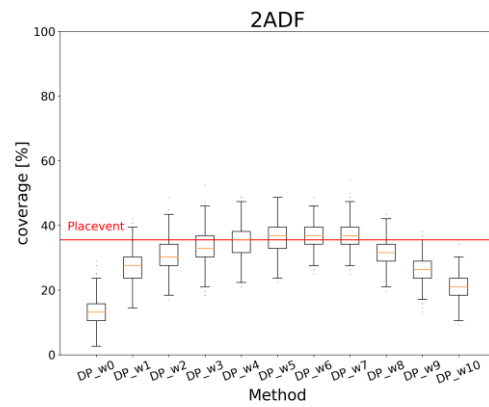
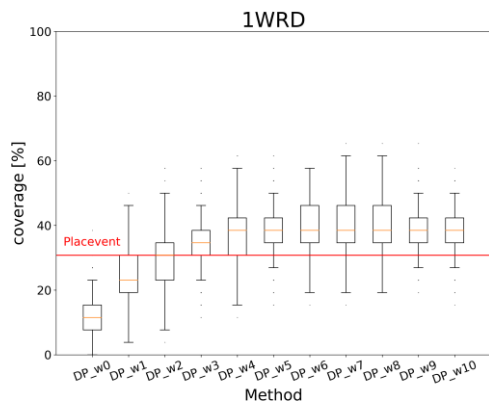
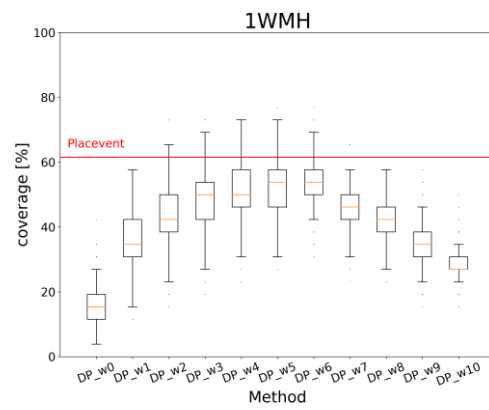
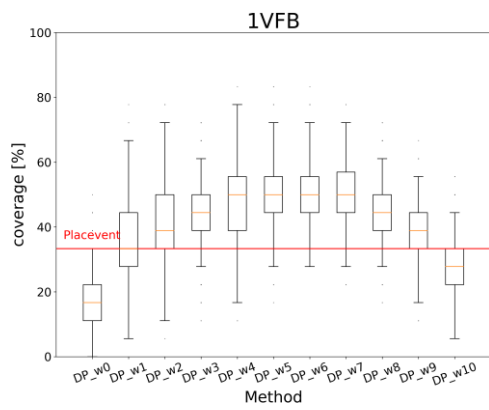
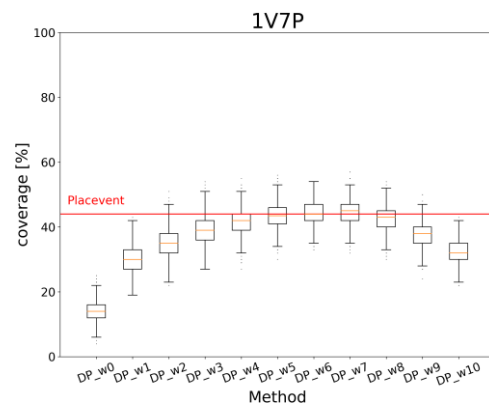
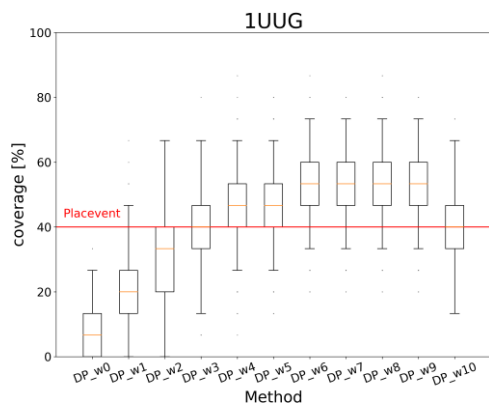
図 S1: ターゲット毎の coverage

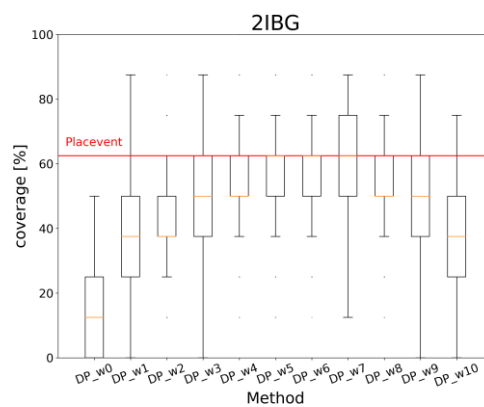
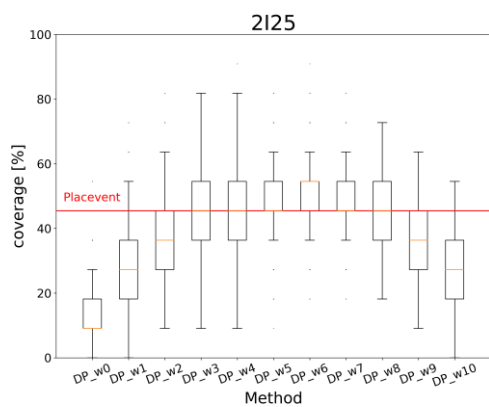
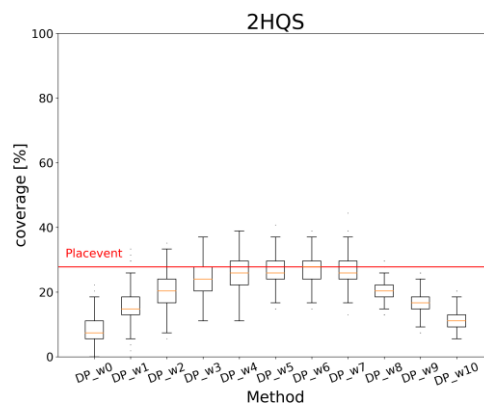
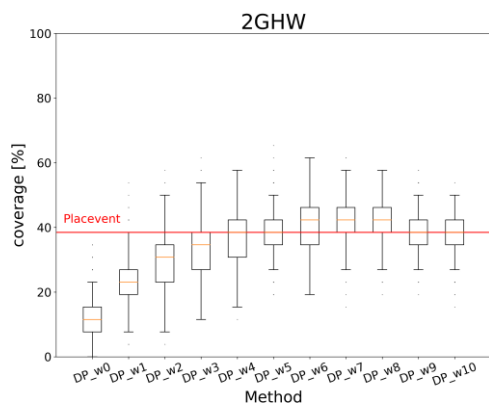
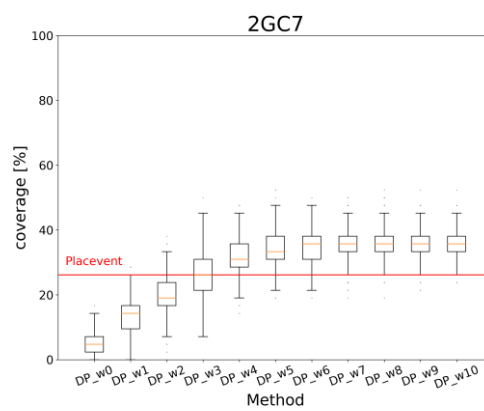
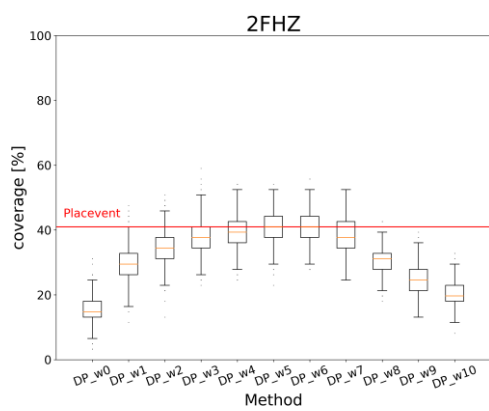
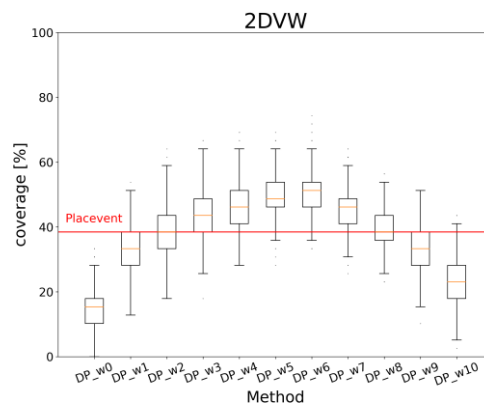
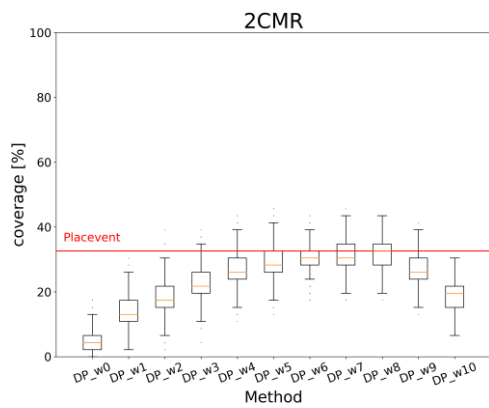
横軸は DP\_w0 から DP\_w10 までの各予測手法を、縦軸は coverage を示す。1000 サンプルの coverage の中央値をオレンジ色で示した。比較のため、各ターゲットに対する Placevent のサンプルの coverage を赤の線で示した。



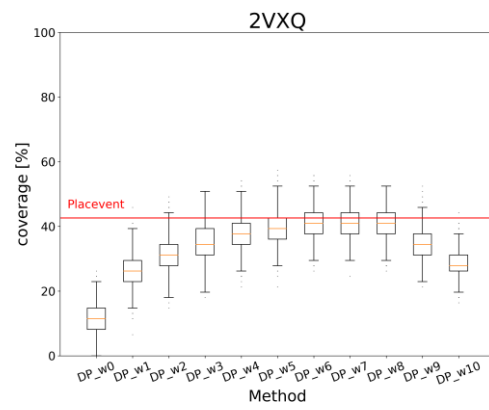
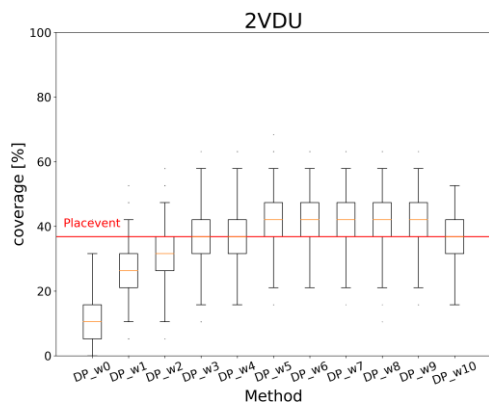
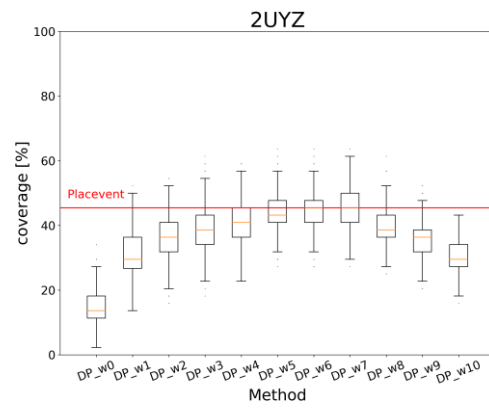
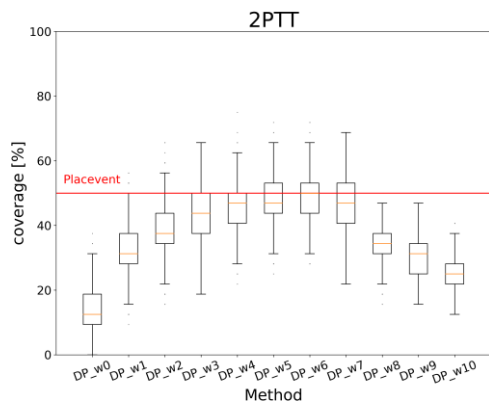
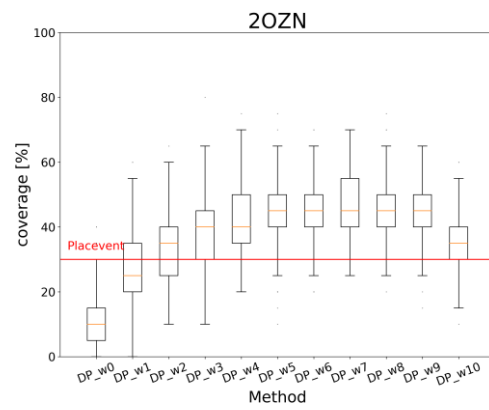
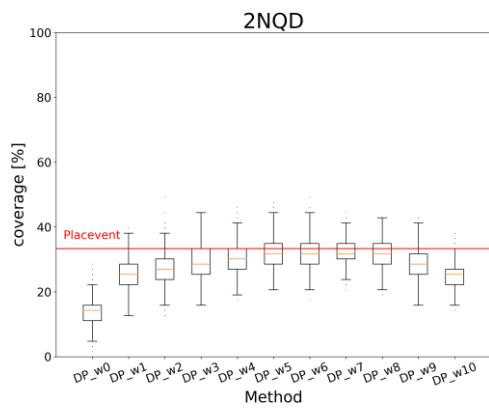
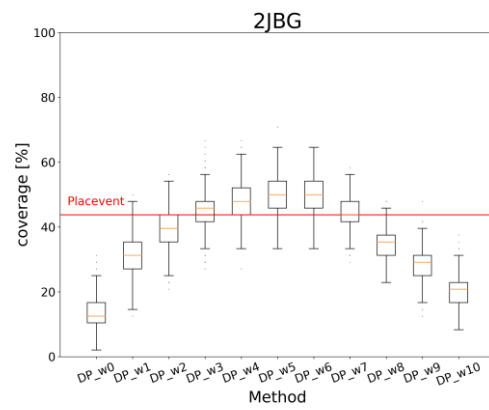
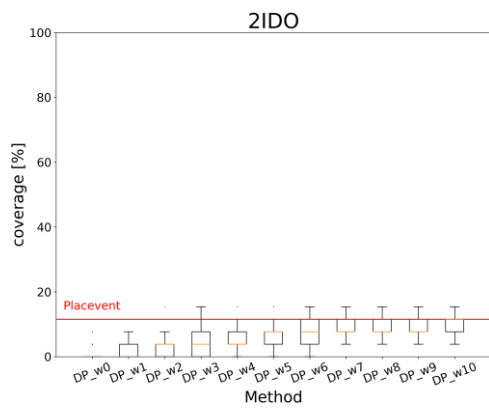


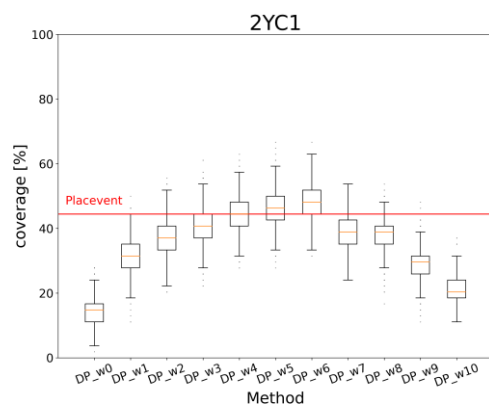
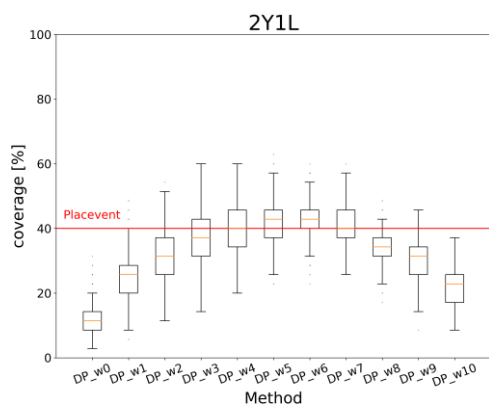
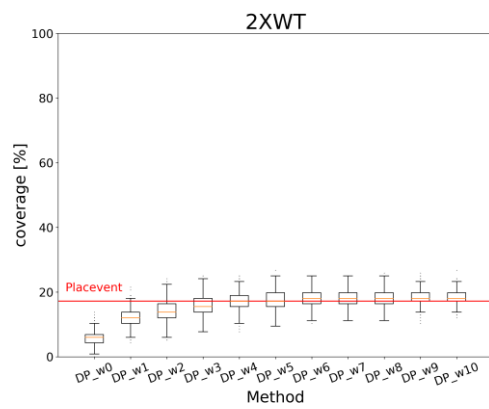
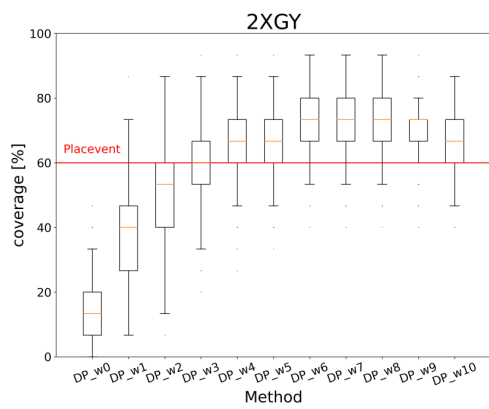
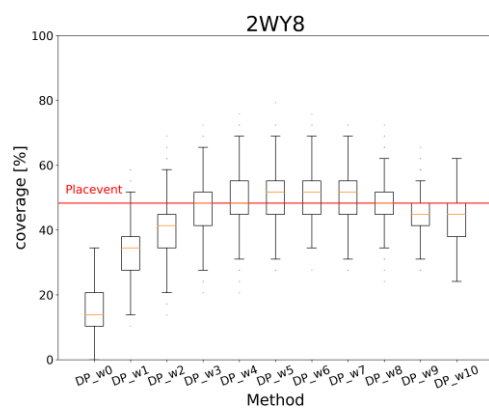
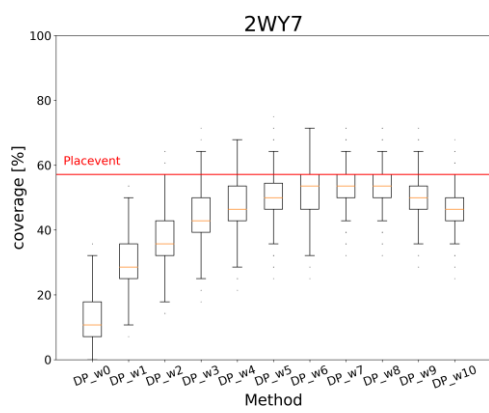
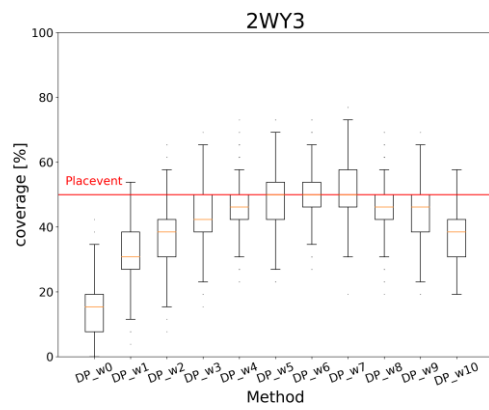
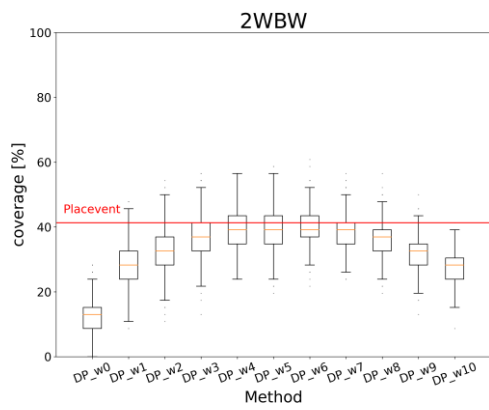


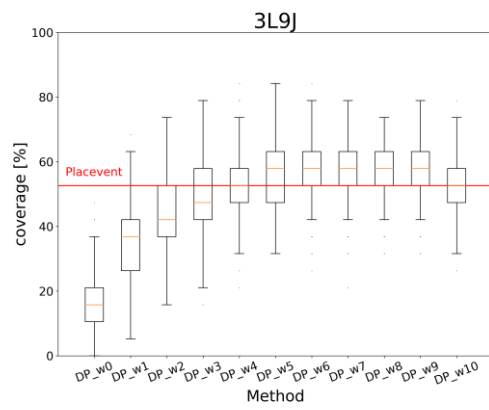
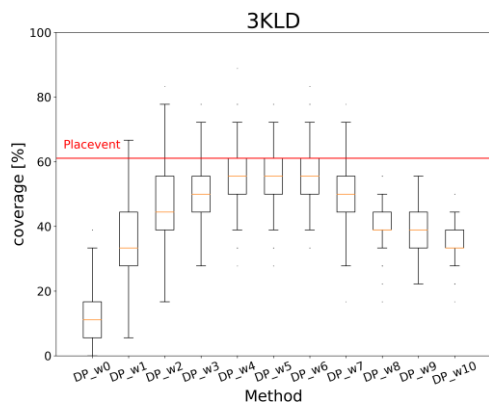
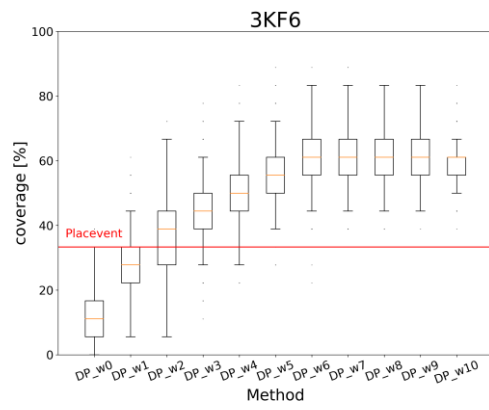
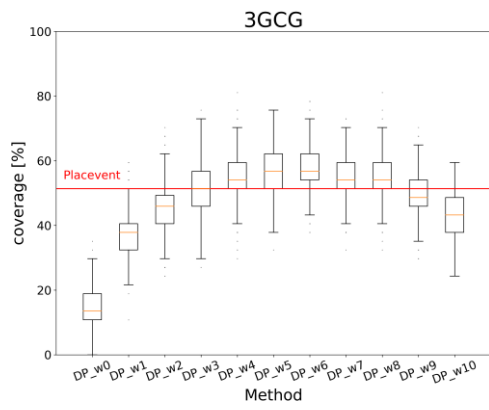
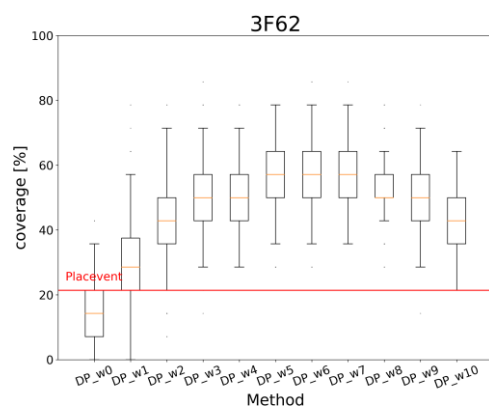
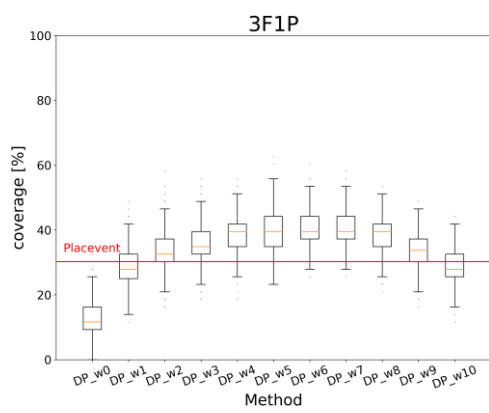
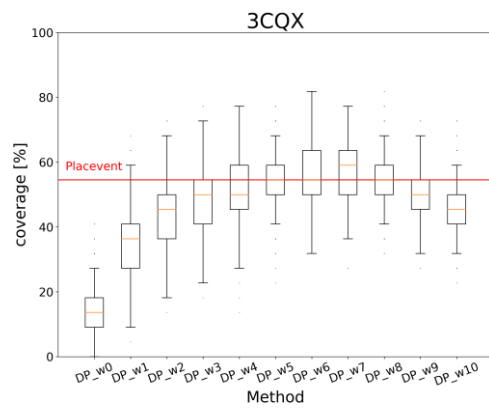
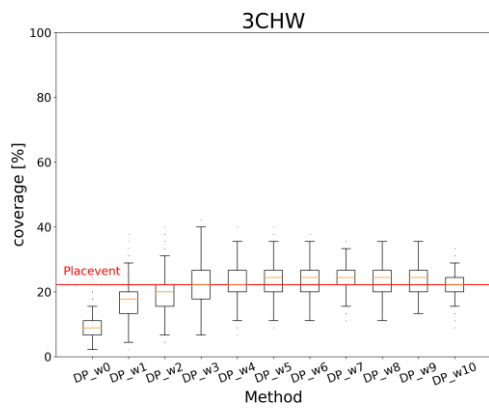


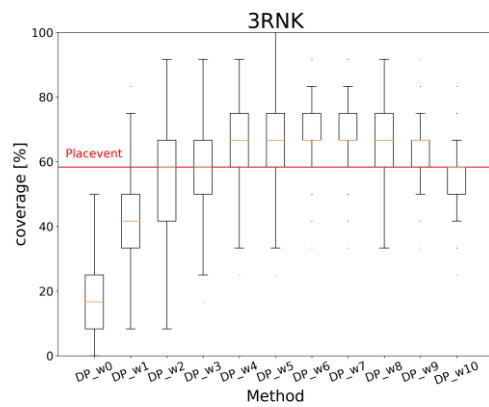
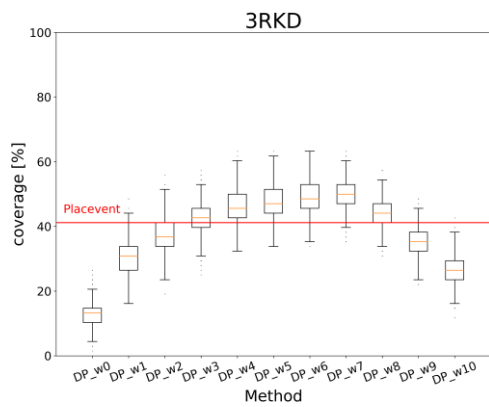
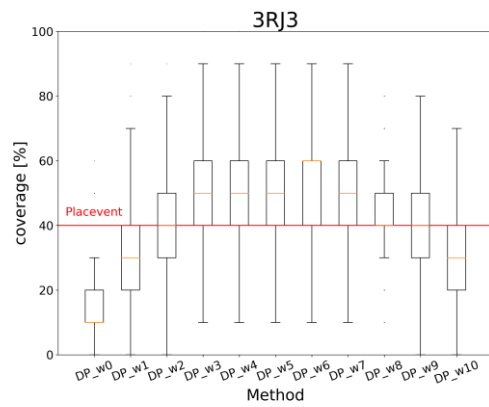
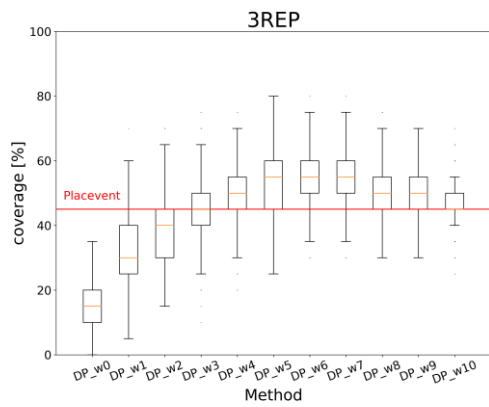
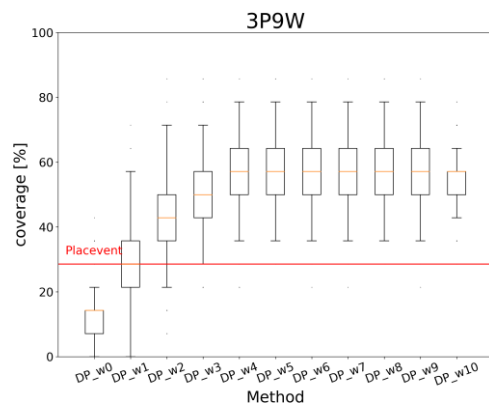
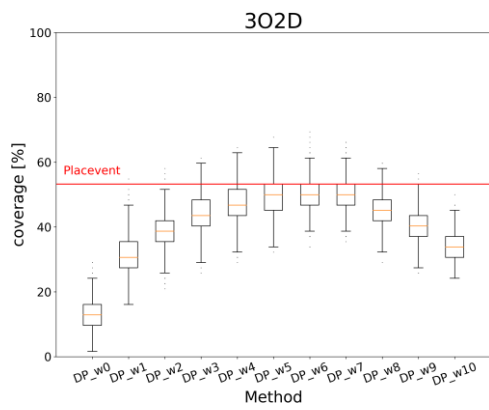
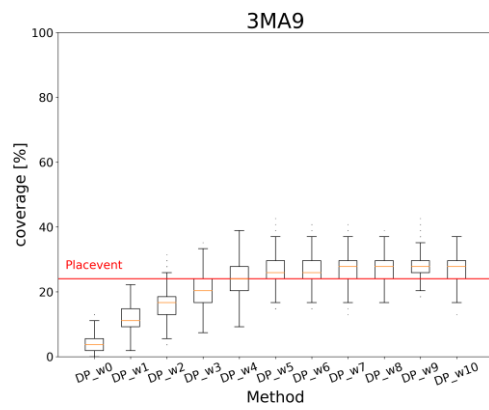
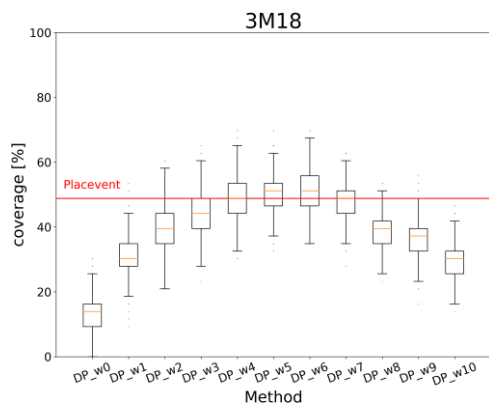


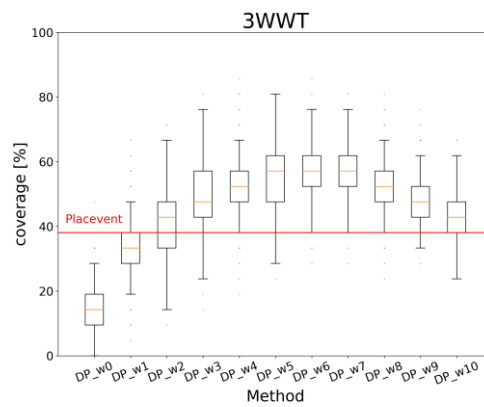
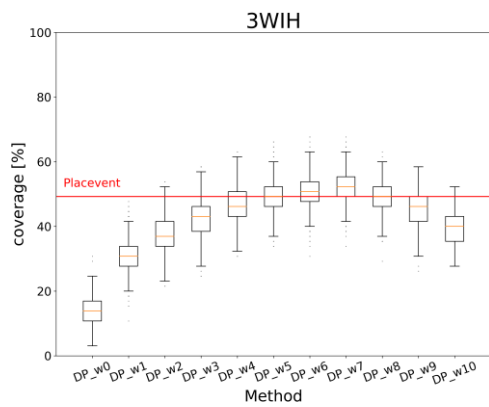
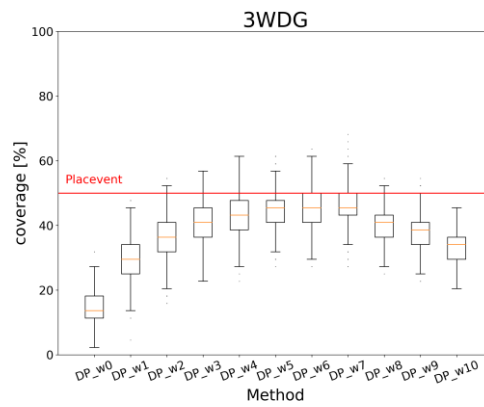
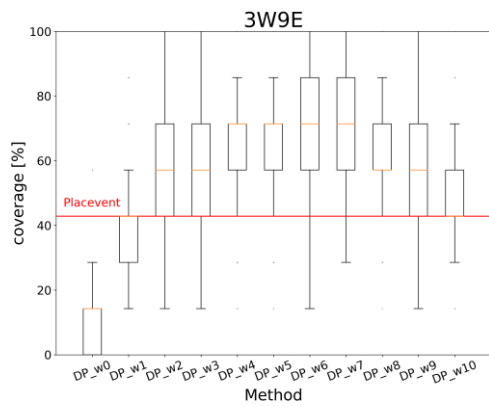
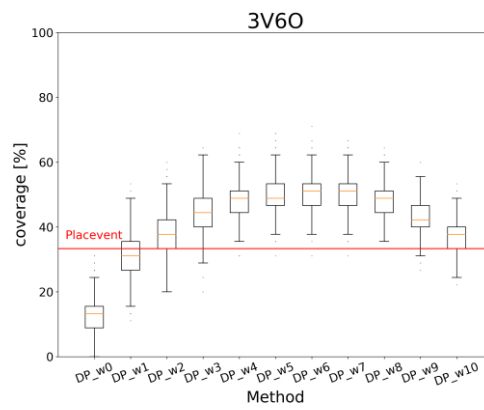
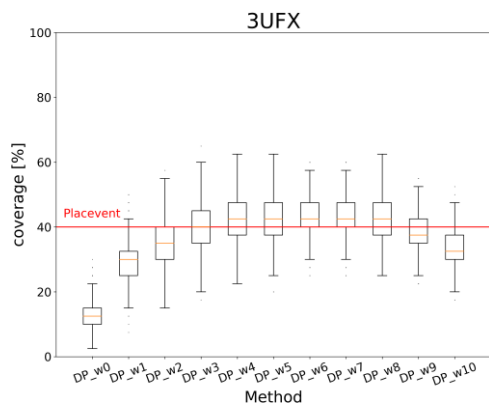
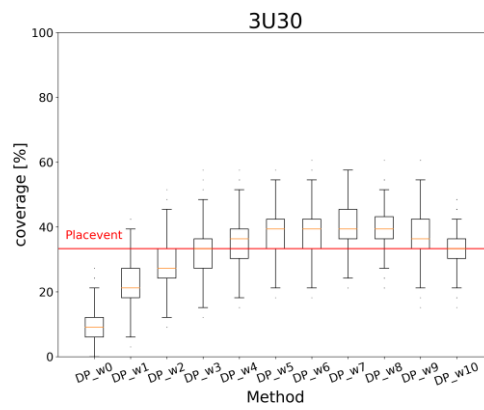
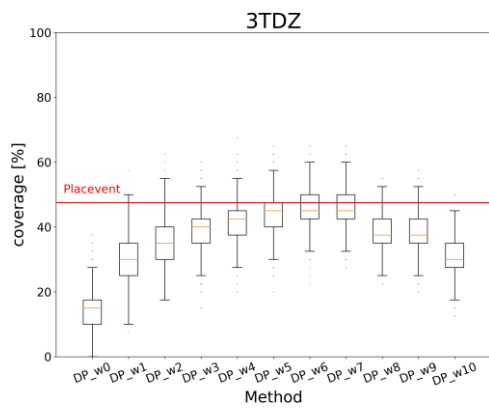


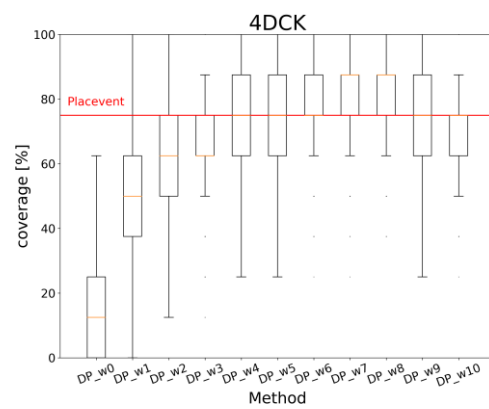
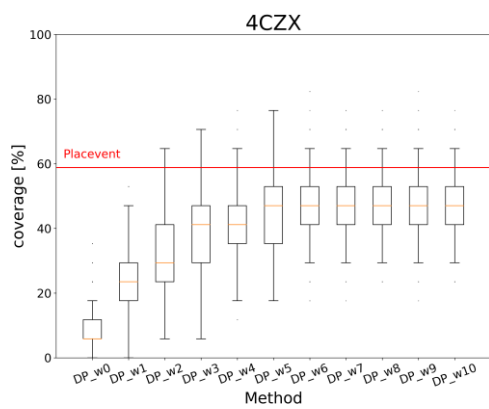
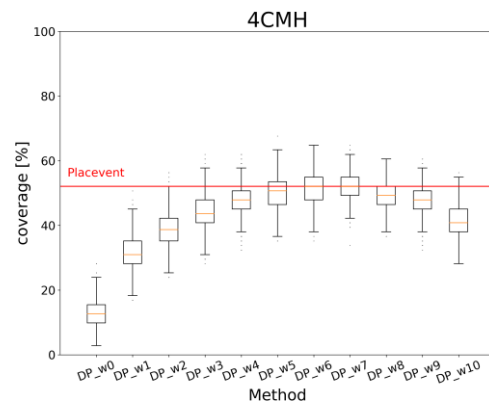
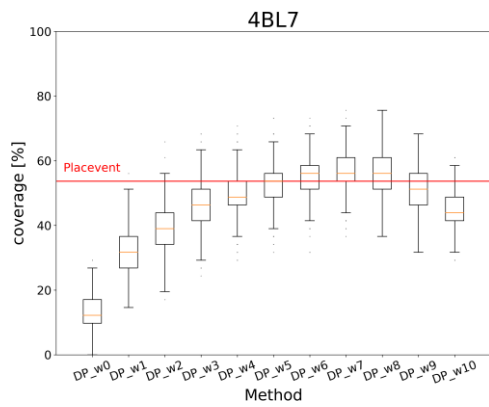
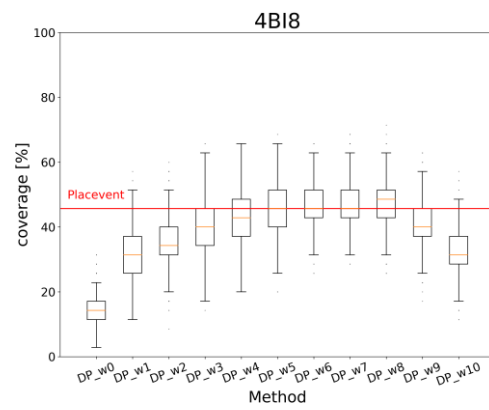
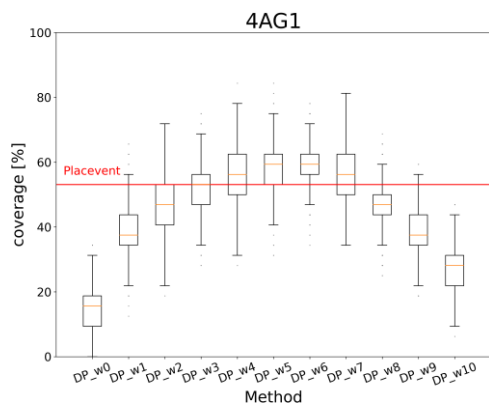
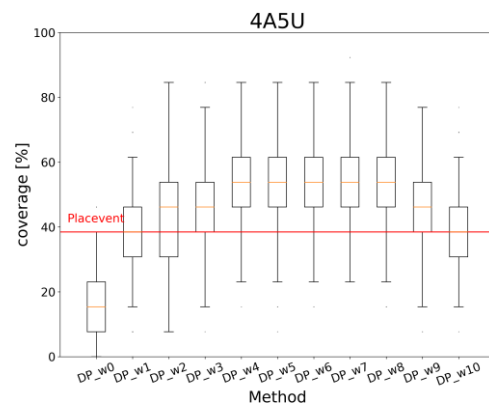
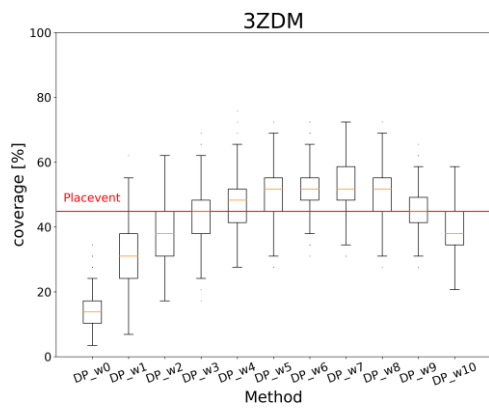


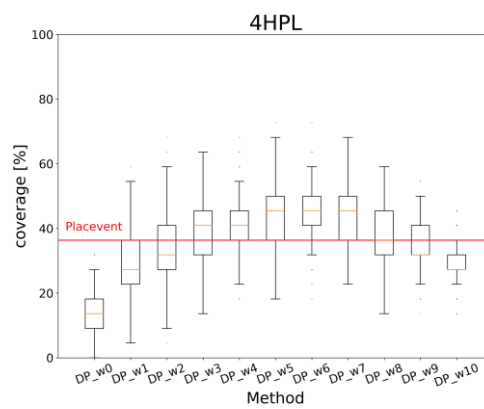
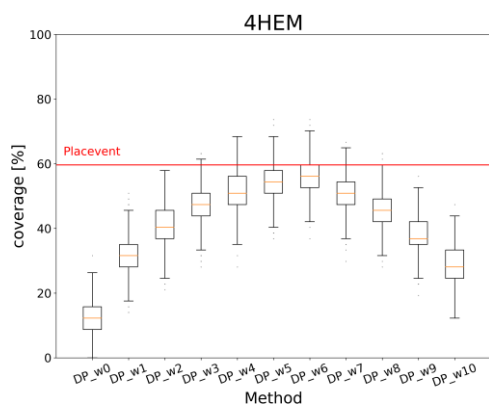
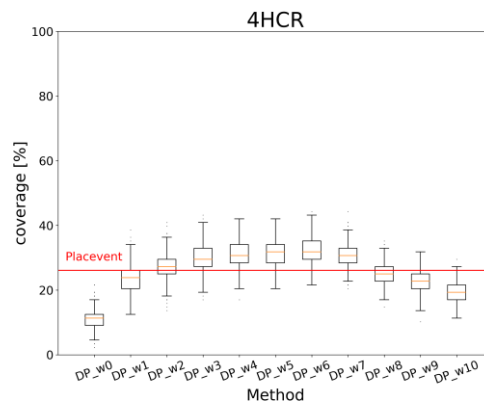
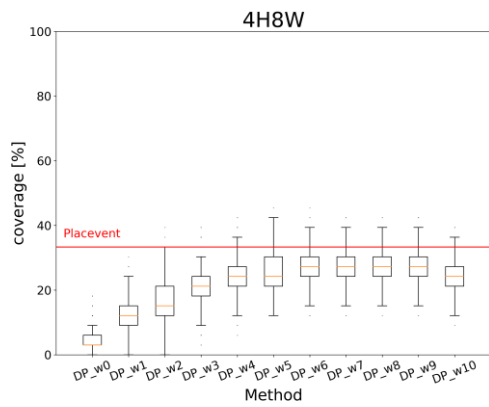
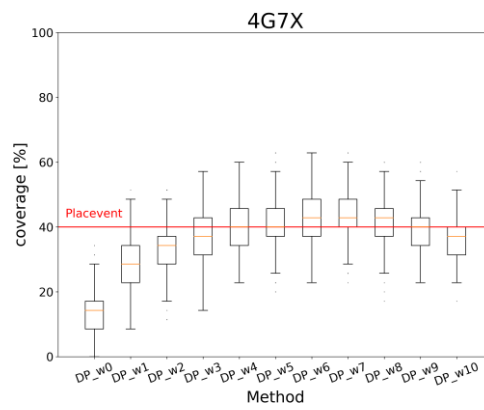
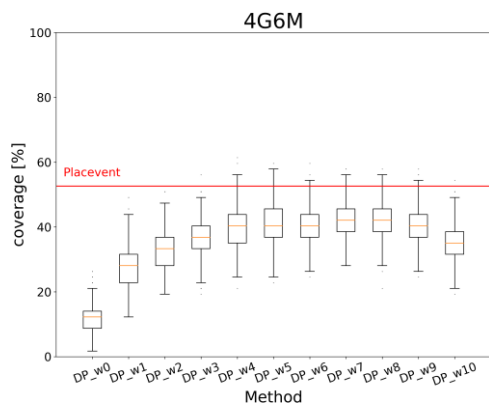
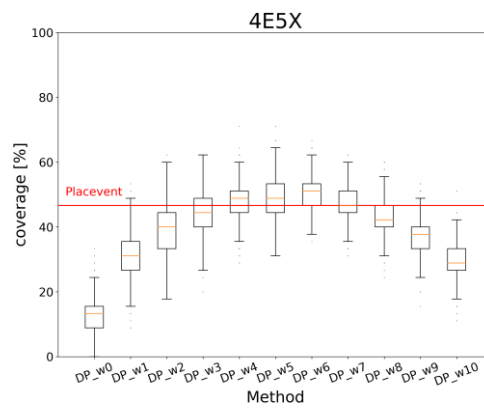
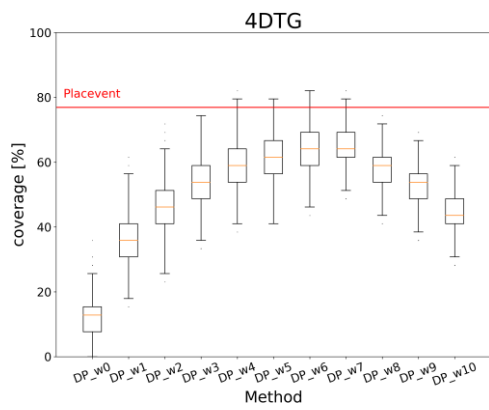


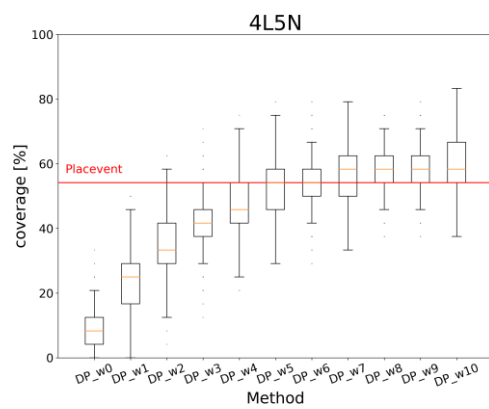
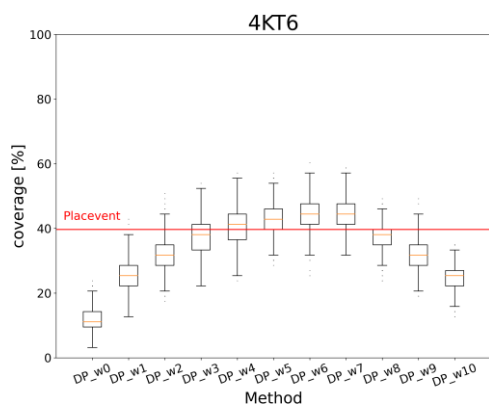
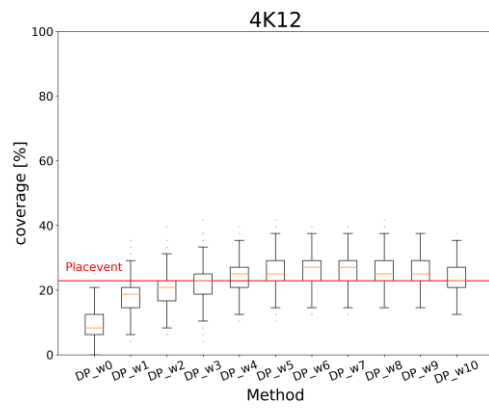
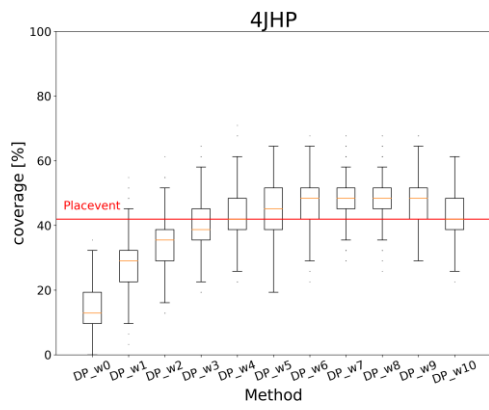
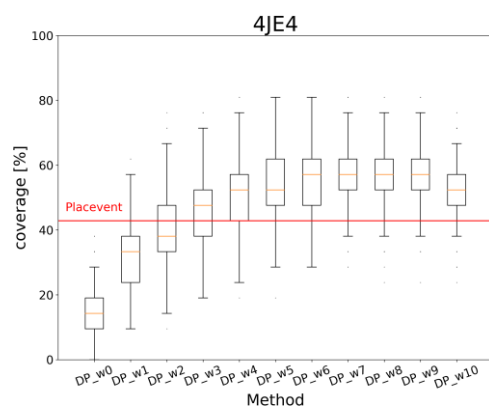
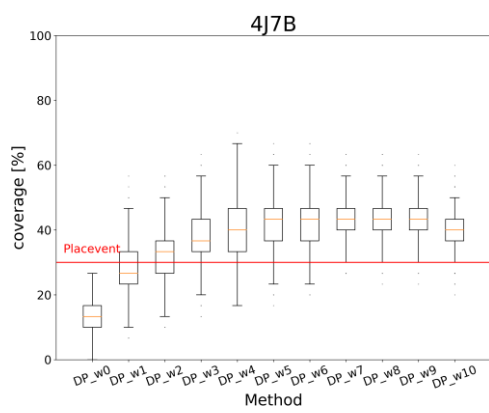
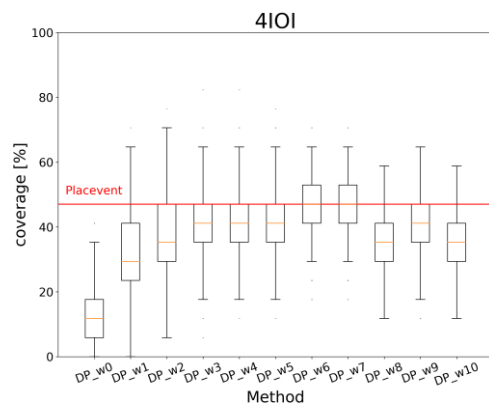
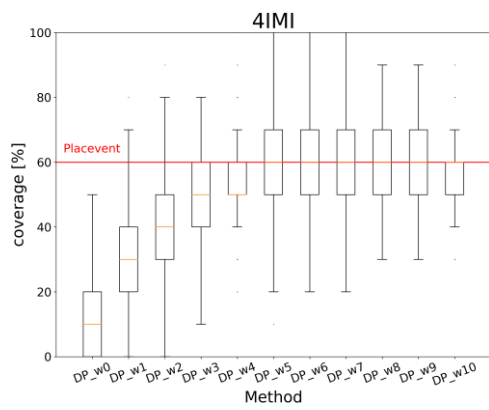




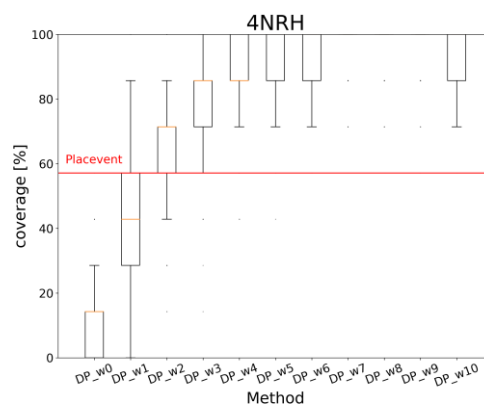
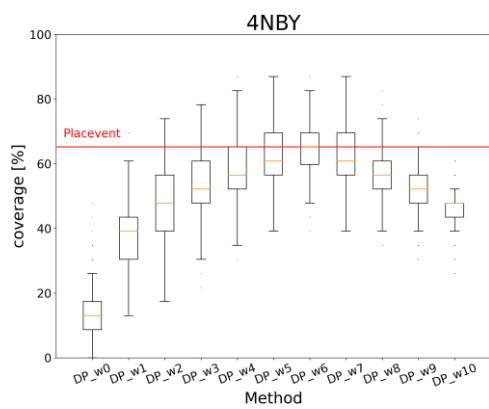
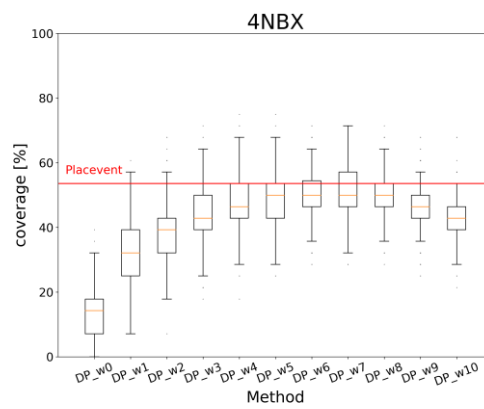
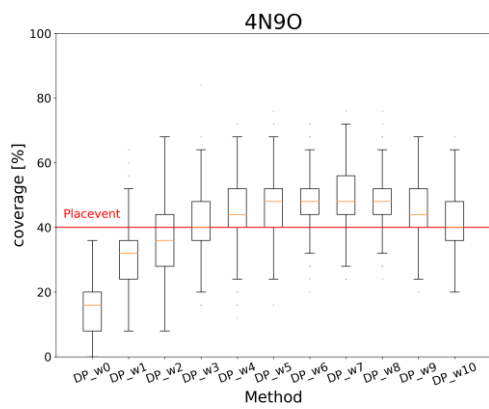
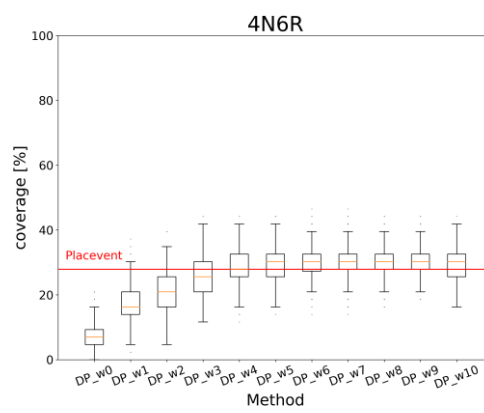
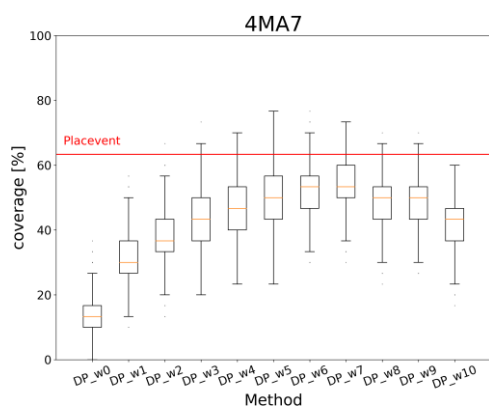
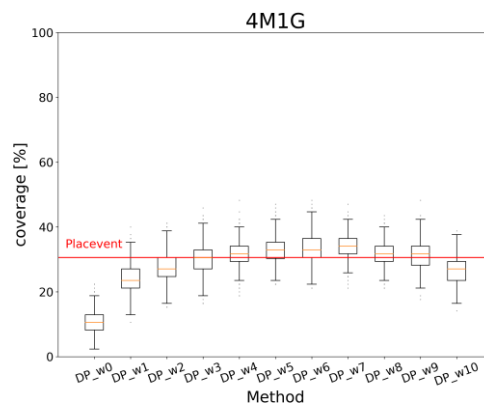
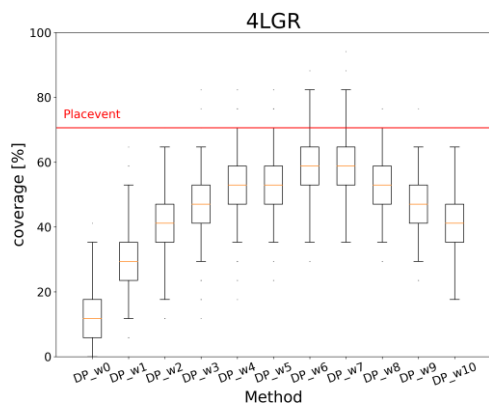


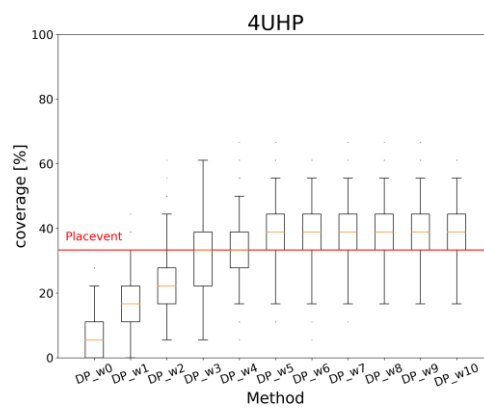
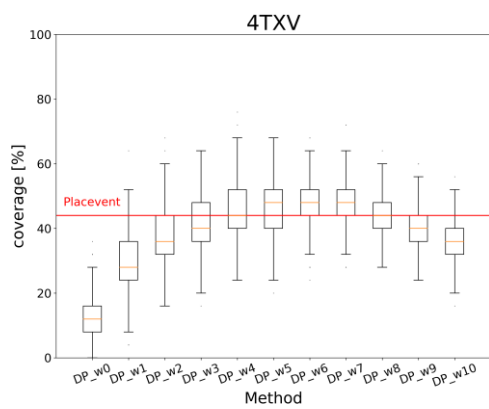
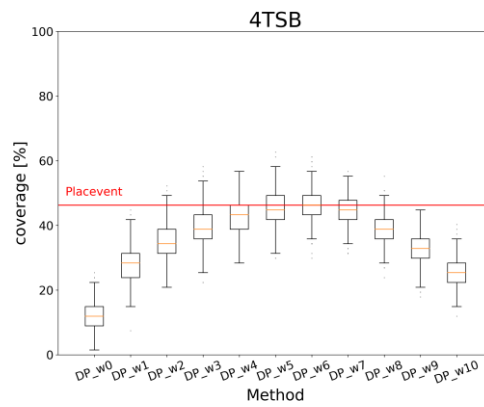
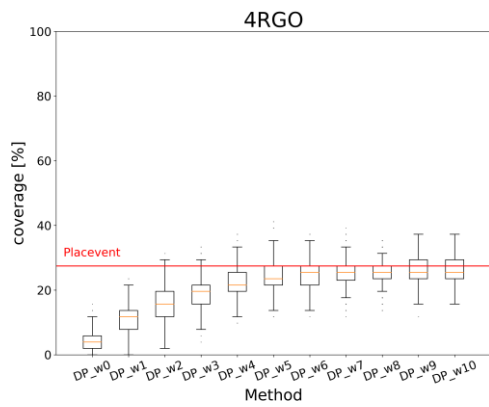
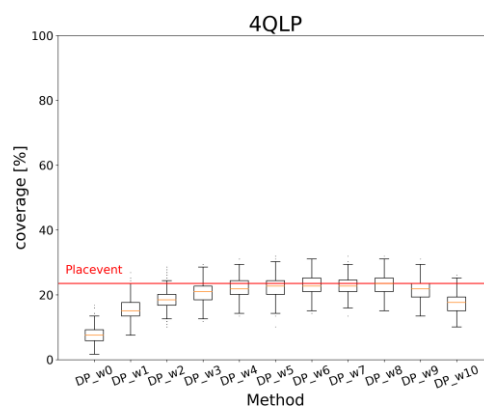
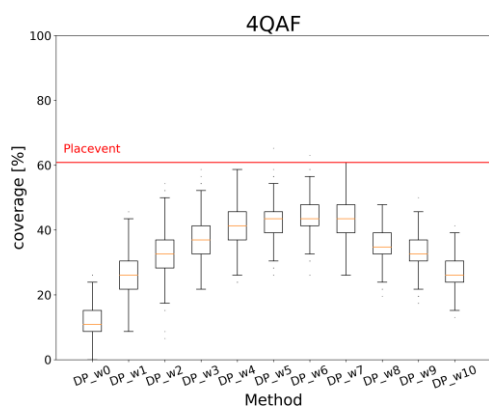
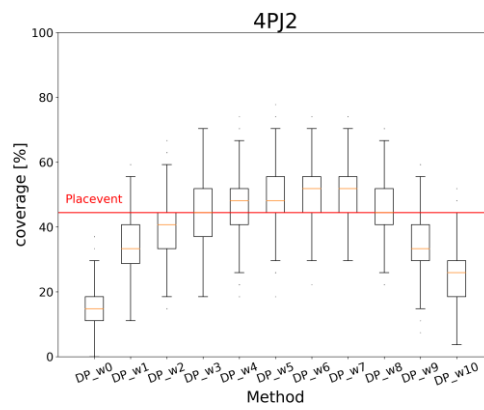
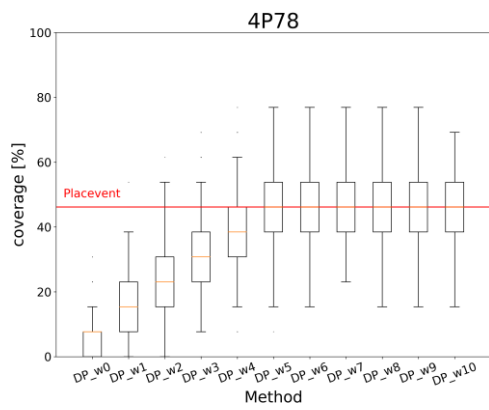


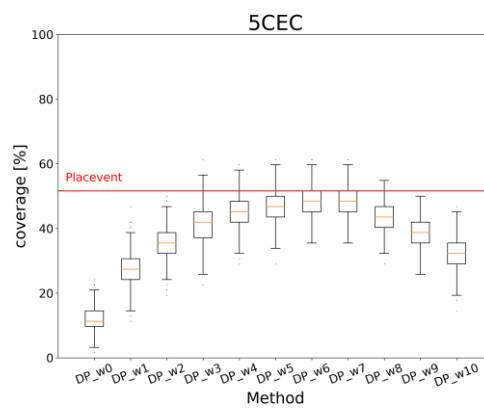
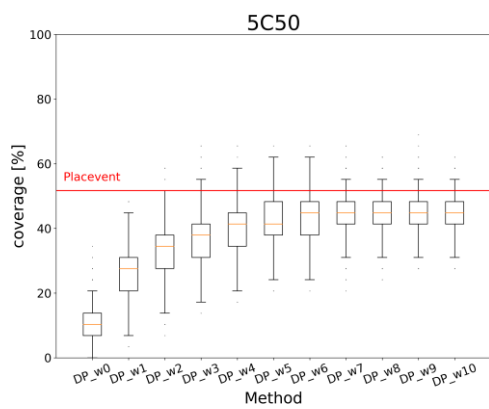
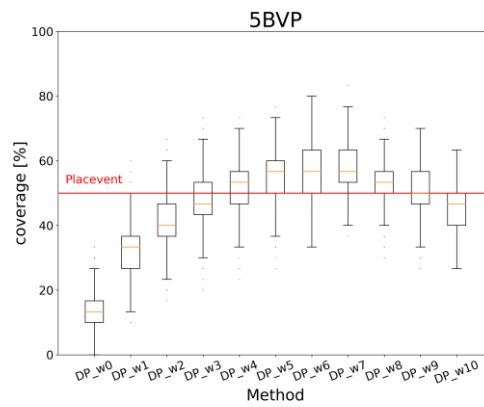
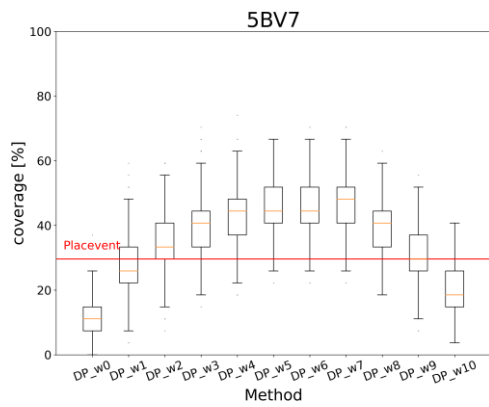
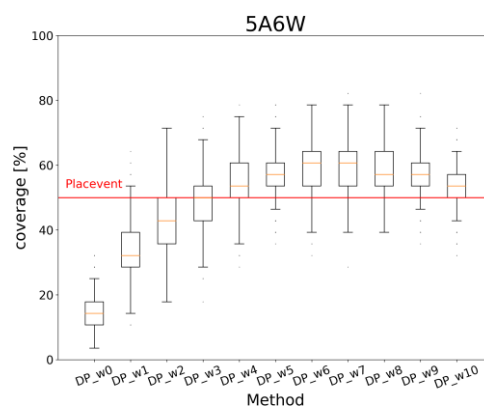
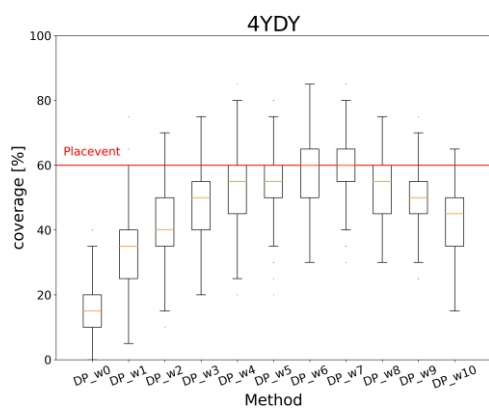
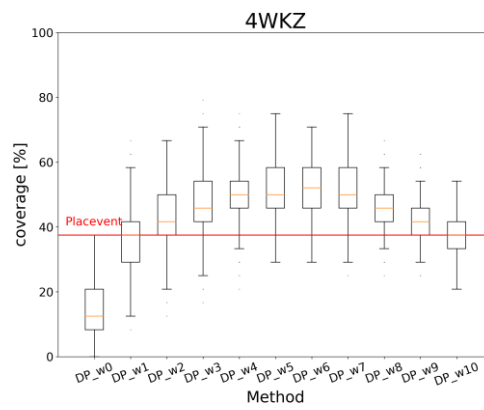
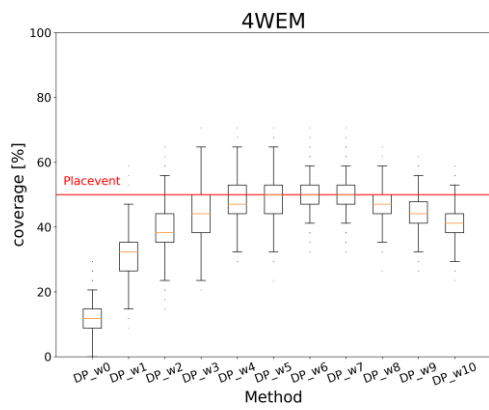


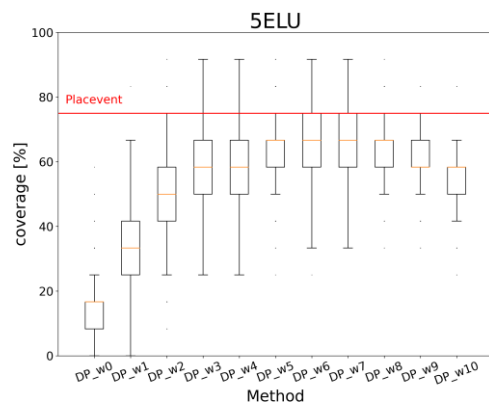
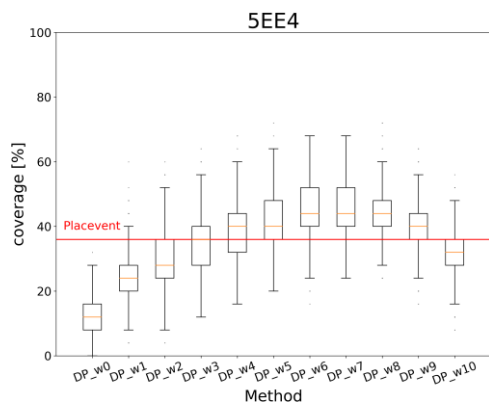
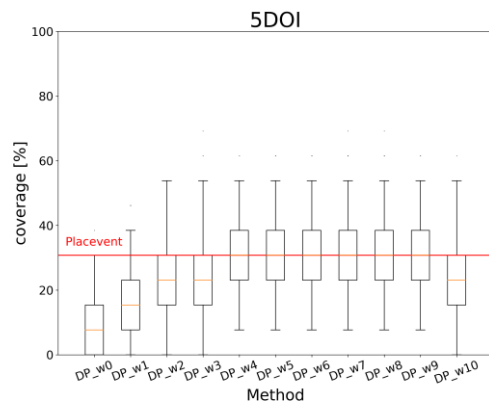
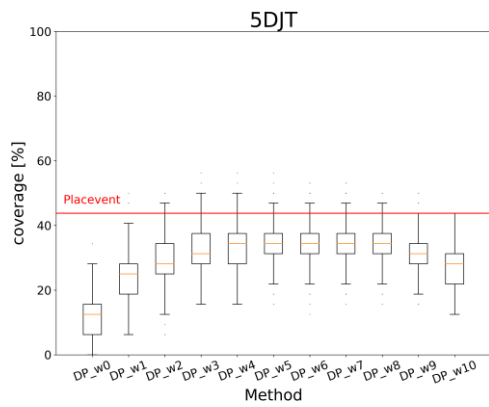
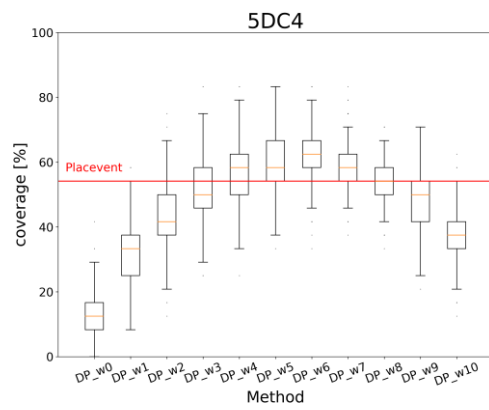
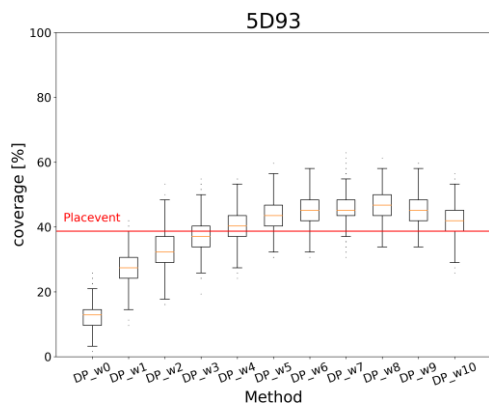
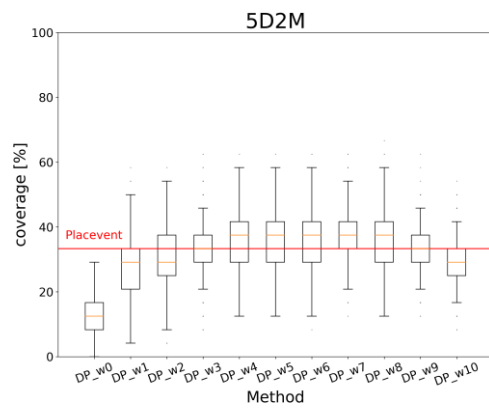
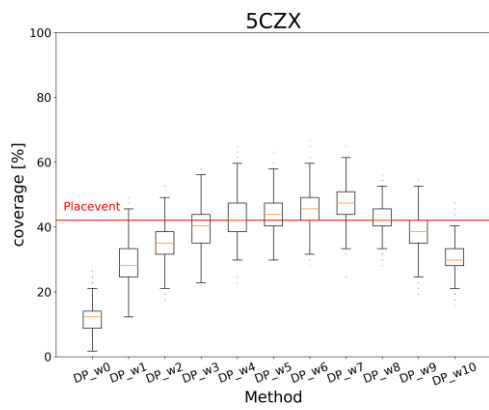


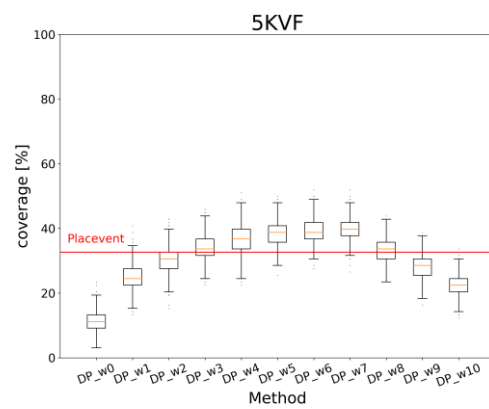
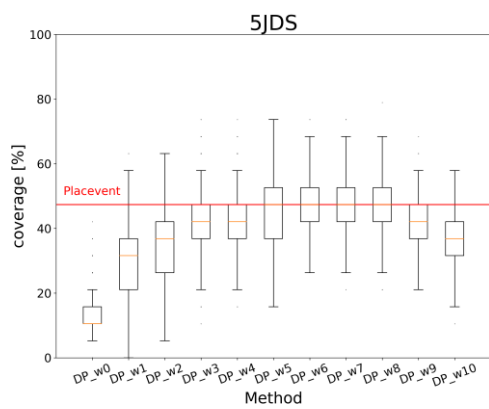
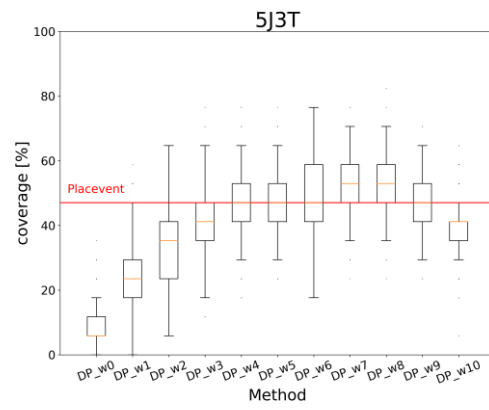
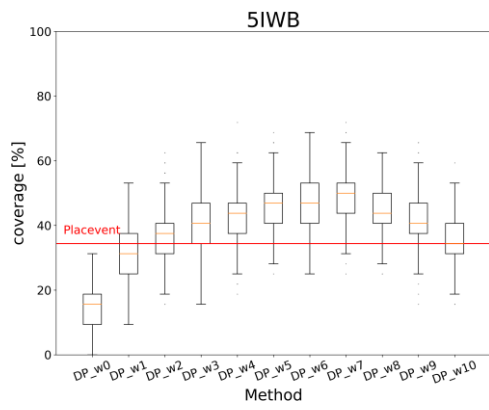
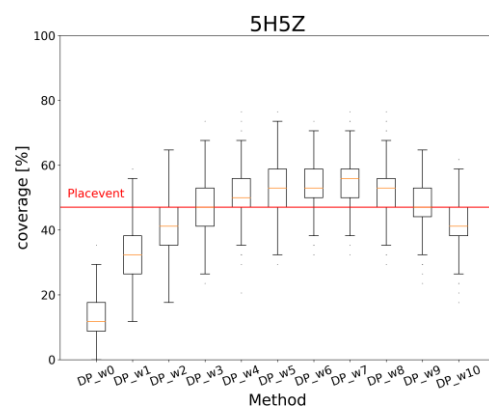
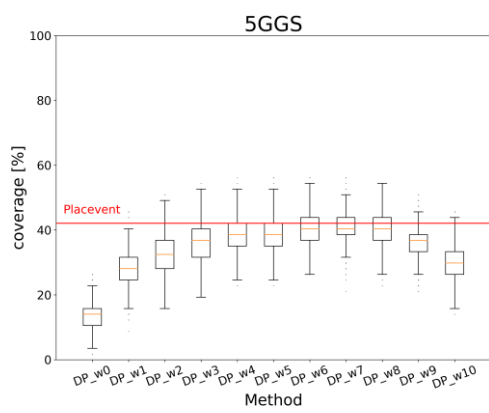
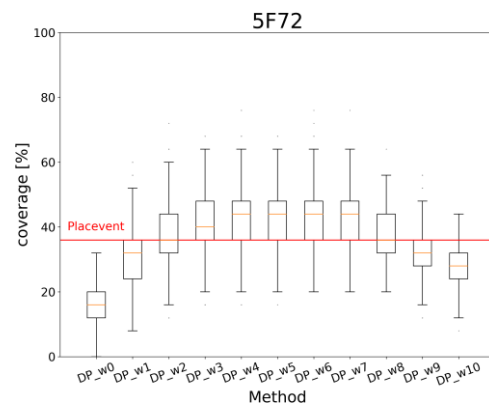
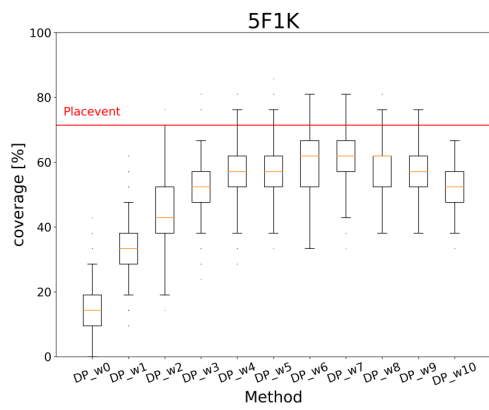


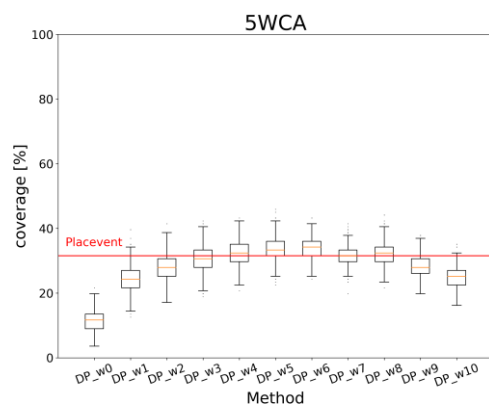
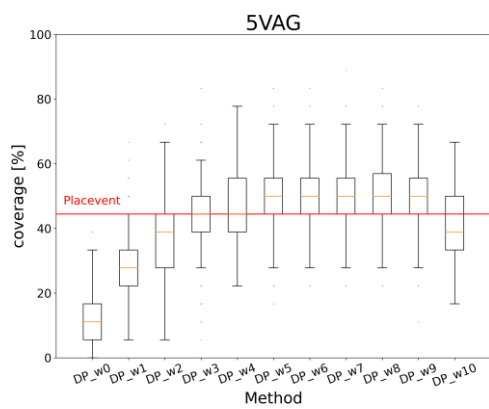
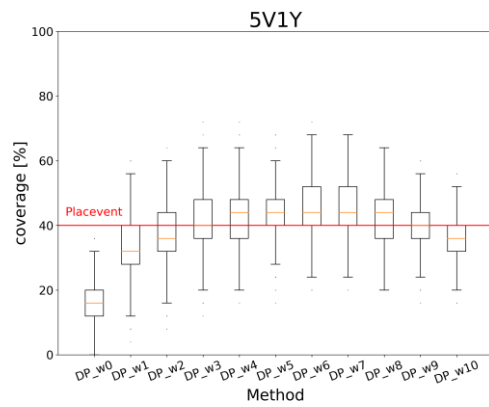
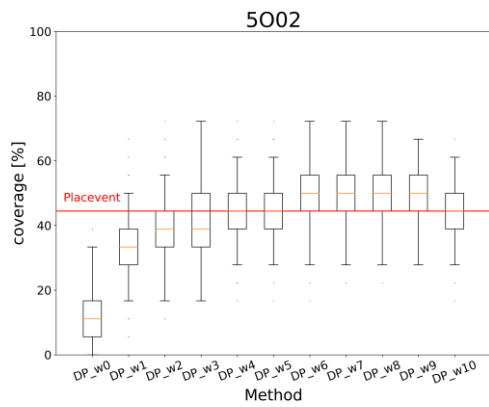
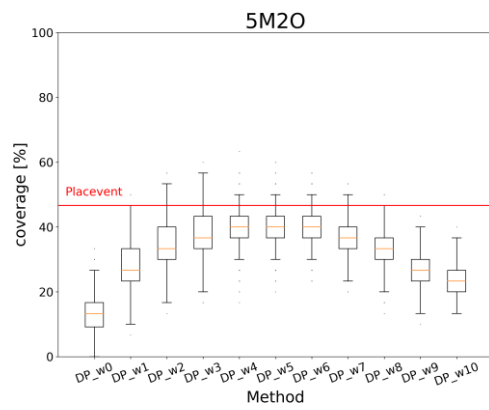
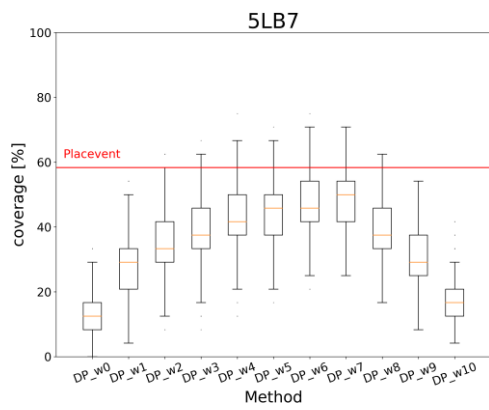
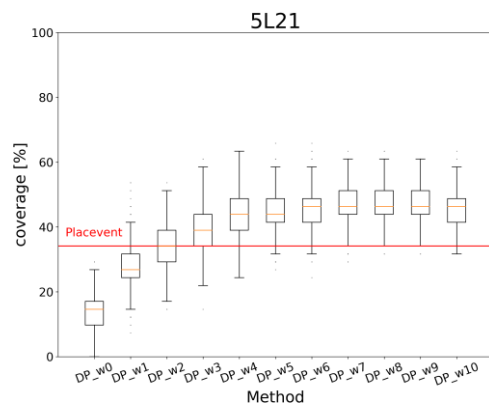
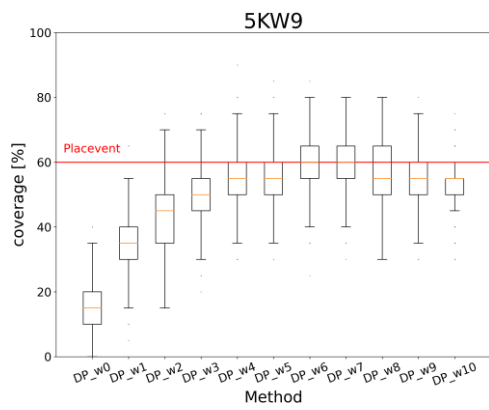


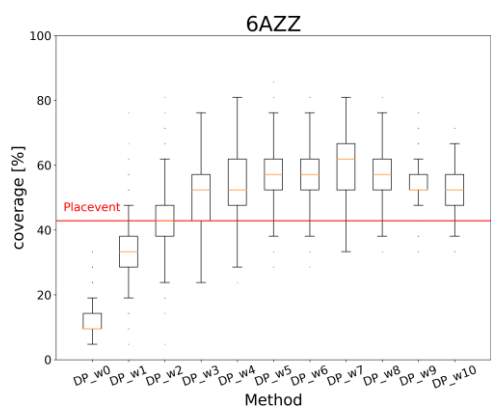
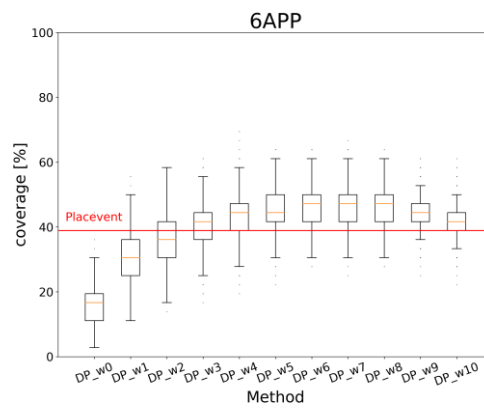
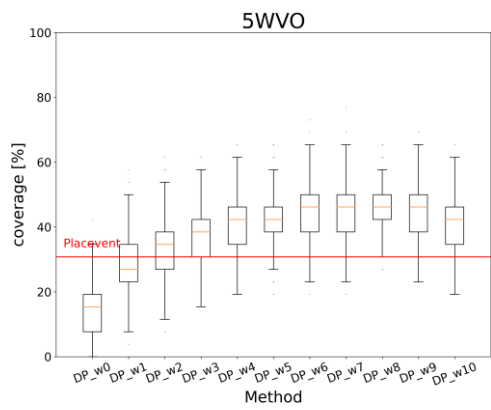












## S-4. Coverage の中央値

サンプル全体の Average coverage の算出に使用した、各ターゲットに対するサンプルの coverage の中央値 (%) を示す (表 S3)。

表 S3: 各手法におけるサンプルの coverage の中央値

各手法によって出力した coverage の中央値 (%) を示す。DroPred (DP\_w0 から DP\_w10) の Average coverage は、1000 サンプルの coverage の中央値となる。Placevent からは各ターゲット 1 サンプルのみが出力されるため、そのサンプルの coverage が coverage の中央値となる。

ID	Placevent	DP_w0	DP_w1	DP_w2	DP_w3	DP_w4	DP_w5	DP_w6	DP_w7	DP_w8	DP_w9	DP_w10
1AY7	64.3	42.9	57.1	64.3	64.3	64.3	64.3	64.3	64.3	64.3	64.3	57.1
1DX5	37.5	37.5	45.8	54.2	54.2	54.2	50.0	50.0	45.8	45.8	41.7	33.3
1DZB	39.1	39.1	47.8	56.5	56.5	56.5	52.2	47.8	43.5	39.1	34.8	30.4
1FYH	44.4	44.4	55.6	66.7	66.7	66.7	66.7	66.7	66.7	66.7	66.7	66.7
1GPQ_1	41.3	37.0	47.8	54.3	54.3	54.3	52.2	47.8	43.5	37.0	30.4	26.1
1GPQ_2	37.8	36.7	47.8	54.4	55.6	55.6	53.3	48.9	43.3	37.8	32.2	26.7
1HQ3	20.0	35.0	52.5	60.0	62.5	65.0	65.0	62.5	57.5	50.0	40.0	30.0
1I2M	36.7	40.0	56.7	63.3	66.7	66.7	66.7	63.3	63.3	60.0	53.3	50.0
1IM3	36.4	36.4	54.5	63.6	63.6	63.6	63.6	63.6	54.5	54.5	54.5	54.5
1IQD	33.3	31.9	40.6	46.4	47.8	47.8	46.4	43.5	39.1	34.8	30.4	23.2
1JTD	31.3	39.6	47.9	54.2	56.3	54.2	52.1	47.9	43.8	39.6	35.4	31.3
1L4D	47.4	36.8	52.6	57.9	57.9	57.9	57.9	52.6	47.4	42.1	42.1	36.8
1LFD	33.3	46.7	53.3	53.3	60.0	53.3	53.3	46.7	46.7	40.0	33.3	33.3
1LPB	58.3	33.3	50.0	58.3	58.3	58.3	58.3	58.3	50.0	41.7	33.3	16.7
1OAK	34.3	34.3	40.0	45.7	45.7	45.7	45.7	42.9	40.0	37.1	31.4	28.6
1OP9	44.4	40.7	55.6	63.0	63.0	63.0	63.0	55.6	51.9	48.1	40.7	37.0
1P2C	42.3	33.8	42.3	47.9	47.9	47.9	45.1	40.8	32.4	23.9	15.5	8.5
1PXV	61.9	38.1	57.1	61.9	66.7	66.7	61.9	61.9	52.4	42.9	33.3	19.0
1TA3	37.5	36.1	44.4	50.0	50.0	50.0	47.2	45.8	43.1	40.3	37.5	34.7
1UOS	57.1	35.7	50.0	57.1	57.1	57.1	57.1	57.1	57.1	57.1	50.0	42.9
1UUG	40.0	40.0	60.0	60.0	66.7	66.7	66.7	66.7	66.7	60.0	53.3	40.0
1V7P	44.0	38.0	50.0	56.0	56.0	56.0	54.0	51.0	47.0	43.0	38.0	32.0



1VFB	33.3	38.9	50.0	55.6	61.1	55.6	55.6	55.6	55.6	44.4	38.9	27.8
1WMH	61.5	42.3	53.8	61.5	65.4	61.5	61.5	57.7	50.0	42.3	34.6	26.9
1WRD	30.8	38.5	42.3	46.2	50.0	50.0	46.2	46.2	42.3	42.3	38.5	38.5
2ADF	35.5	36.8	47.4	53.9	53.9	52.6	50.0	44.7	39.5	31.6	26.3	21.1
2BKY_1	44.8	37.9	48.3	55.2	55.2	55.2	51.7	48.3	44.8	37.9	24.1	17.2
2BKY_2	42.4	35.6	47.5	52.5	52.5	52.5	50.8	45.8	40.7	33.9	25.4	18.6
2CMR	32.6	39.1	52.2	58.7	58.7	58.7	56.5	52.2	43.5	34.8	26.1	19.6
2DVW	38.5	38.5	51.3	59.0	59.0	59.0	59.0	53.8	48.7	41.0	33.3	23.1
2FHZ	41.0	39.3	47.5	52.5	52.5	52.5	49.2	44.3	37.7	31.1	24.6	19.7
2GC7	26.2	40.5	57.1	61.9	61.9	61.9	59.5	57.1	54.8	50.0	45.2	40.5
2GHW	38.5	38.5	46.2	53.8	53.8	53.8	50.0	46.2	46.2	42.3	38.5	38.5
2HQS	27.8	38.9	44.4	51.9	51.9	51.9	48.1	42.6	31.5	24.1	16.7	11.1
2I25	45.5	36.4	54.5	54.5	54.5	54.5	54.5	54.5	45.5	45.5	36.4	27.3
2IBG	62.5	37.5	50.0	62.5	62.5	62.5	62.5	62.5	62.5	62.5	50.0	37.5
2ID0	11.5	38.5	50.0	57.7	57.7	57.7	53.8	50.0	46.2	42.3	34.6	30.8
2JBG	43.8	39.6	52.1	58.3	60.4	58.3	56.3	52.1	45.8	37.5	29.2	20.8
2NQD	33.3	36.5	44.4	50.8	49.2	47.6	44.4	41.3	38.1	33.3	28.6	25.4
2OZN	30.0	40.0	50.0	55.0	55.0	60.0	55.0	55.0	50.0	50.0	45.0	35.0
2PTT	50.0	37.5	53.1	59.4	62.5	62.5	59.4	56.3	46.9	37.5	31.3	25.0
2UYZ	45.5	36.4	47.7	54.5	56.8	56.8	54.5	52.3	50.0	43.2	36.4	29.5
2VDU	36.8	36.8	42.1	47.4	52.6	47.4	47.4	47.4	42.1	42.1	42.1	36.8
2VXQ	42.6	36.1	47.5	54.1	55.7	54.1	54.1	50.8	45.9	41.0	34.4	27.9
2WBW	41.3	37.0	47.8	54.3	54.3	52.2	50.0	45.7	41.3	37.0	32.6	28.3
2WY3	50.0	38.5	50.0	57.7	57.7	57.7	57.7	57.7	53.8	50.0	46.2	38.5
2WY7	57.1	39.3	50.0	57.1	60.7	60.7	60.7	57.1	57.1	53.6	50.0	46.4
2WY8	48.3	37.9	51.7	58.6	58.6	58.6	58.6	55.2	51.7	48.3	44.8	44.8
2XGY	60.0	40.0	60.0	66.7	73.3	73.3	73.3	73.3	73.3	73.3	73.3	66.7
2XWT	17.2	36.2	45.7	51.7	51.7	50.0	46.6	42.2	37.1	31.9	26.7	22.4
2Y1L	40.0	40.0	48.6	54.3	54.3	54.3	51.4	48.6	42.9	37.1	31.4	22.9
2YC1	44.4	37.0	51.9	57.4	59.3	59.3	57.4	55.6	48.1	38.9	29.6	20.4
3CHW	22.2	40.0	48.9	55.6	53.3	51.1	46.7	42.2	37.8	31.1	26.7	24.4
3CQX	54.5	45.5	54.5	59.1	63.6	63.6	63.6	59.1	59.1	59.1	50.0	45.5
3F1P	30.2	34.9	46.5	53.5	51.2	51.2	48.8	44.2	41.9	39.5	32.6	27.9
3F62	21.4	42.9	57.1	64.3	64.3	64.3	64.3	64.3	57.1	57.1	50.0	42.9

3GCG	51.4	40.5	56.8	62.2	64.9	64.9	64.9	62.2	59.5	54.1	48.6	43.2
3KF6	33.3	38.9	55.6	66.7	66.7	66.7	66.7	66.7	66.7	66.7	61.1	61.1
3KLD	61.1	44.4	61.1	66.7	66.7	72.2	66.7	61.1	50.0	44.4	38.9	33.3
3L9J	52.6	36.8	52.6	57.9	63.2	63.2	63.2	63.2	63.2	57.9	57.9	52.6
3M18	48.8	39.5	51.2	58.1	60.5	60.5	58.1	55.8	48.8	44.2	37.2	30.2
3MA9	24.1	35.2	48.1	55.6	57.4	57.4	57.4	55.6	51.9	44.4	37.0	31.5
3O2D	53.2	37.1	50.0	58.1	58.1	59.7	58.1	54.8	51.6	46.8	40.3	33.9
3P9W	28.6	42.9	50.0	57.1	57.1	57.1	57.1	57.1	57.1	57.1	57.1	57.1
3REP	45.0	40.0	50.0	60.0	60.0	60.0	60.0	60.0	55.0	55.0	50.0	45.0
3RJ3	40.0	40.0	50.0	60.0	60.0	60.0	60.0	60.0	50.0	50.0	40.0	30.0
3RKD	41.2	39.7	52.9	60.3	61.8	61.8	58.8	55.9	51.5	44.1	35.3	26.5
3RNK	58.3	41.7	66.7	66.7	75.0	75.0	75.0	75.0	75.0	66.7	66.7	58.3
3TDZ	47.5	37.5	50.0	57.5	57.5	57.5	55.0	55.0	50.0	45.0	37.5	30.0
3U30	33.3	33.3	42.4	48.5	48.5	48.5	48.5	48.5	45.5	42.4	36.4	33.3
3UFX	40.0	40.0	50.0	55.0	57.5	57.5	55.0	52.5	47.5	42.5	37.5	32.5
3V60	33.3	37.8	53.3	60.0	60.0	60.0	60.0	55.6	51.1	48.9	42.2	37.8
3W9E	42.9	42.9	57.1	71.4	71.4	71.4	71.4	71.4	71.4	57.1	57.1	42.9
3WDG	50.0	38.6	47.7	54.5	56.8	56.8	54.5	52.3	50.0	43.2	38.6	34.1
3WIH	49.2	38.5	50.8	56.9	58.5	60.0	58.5	56.9	53.8	50.8	46.2	40.0
3WWT	38.1	38.1	52.4	61.9	66.7	66.7	66.7	61.9	57.1	52.4	47.6	42.9
3ZDM	44.8	37.9	51.7	58.6	62.1	62.1	62.1	58.6	55.2	51.7	44.8	37.9
4A5U	38.5	38.5	53.8	61.5	61.5	61.5	61.5	53.8	53.8	53.8	46.2	38.5
4AG1	53.1	40.6	56.3	62.5	62.5	62.5	65.6	62.5	56.3	46.9	37.5	28.1
4BI8	45.7	37.1	51.4	60.0	60.0	60.0	60.0	60.0	54.3	48.6	40.0	31.4
4BL7	53.7	39.0	56.1	61.0	63.4	63.4	63.4	63.4	61.0	56.1	51.2	43.9
4CMH	52.1	39.4	53.5	62.0	62.0	63.4	62.0	59.2	56.3	53.5	47.9	40.8
4CZX	58.8	41.2	52.9	58.8	58.8	64.7	58.8	58.8	58.8	58.8	52.9	47.1
4DCK	75.0	50.0	62.5	75.0	75.0	87.5	87.5	87.5	87.5	87.5	75.0	75.0
4DTG	76.9	41.0	56.4	66.7	69.2	69.2	69.2	69.2	64.1	59.0	53.8	43.6
4E5X	46.7	40.0	53.3	60.0	60.0	62.2	60.0	55.6	51.1	44.4	37.8	28.9
4G6M	52.6	38.6	52.6	57.9	59.6	59.6	57.9	54.4	50.9	45.6	40.4	35.1
4G7X	40.0	37.1	48.6	54.3	54.3	54.3	51.4	51.4	45.7	42.9	40.0	37.1
4H8W	33.3	39.4	51.5	60.6	63.6	63.6	60.6	57.6	48.5	42.4	33.3	24.2
4HCR	26.1	28.4	36.4	40.9	40.9	39.8	38.6	35.2	31.8	27.3	22.7	19.3

4HEM	59.6	38.6	54.4	61.4	63.2	63.2	63.2	59.6	52.6	45.6	36.8	28.1
4HPL	36.4	36.4	50.0	54.5	59.1	59.1	59.1	54.5	50.0	40.9	31.8	27.3
4IMI	60.0	50.0	60.0	70.0	70.0	70.0	70.0	70.0	70.0	60.0	60.0	60.0
4IOI	47.1	41.2	52.9	58.8	58.8	58.8	52.9	52.9	47.1	41.2	41.2	35.3
4J7B	30.0	36.7	50.0	56.7	56.7	56.7	56.7	53.3	50.0	46.7	43.3	40.0
4JE4	42.9	38.1	52.4	57.1	61.9	61.9	61.9	61.9	61.9	61.9	57.1	52.4
4JHP	41.9	38.7	51.6	58.1	61.3	61.3	61.3	61.3	58.1	51.6	48.4	41.9
4K12	22.9	35.4	45.8	52.1	52.1	52.1	50.0	45.8	41.7	35.4	29.2	22.9
4KT6	39.7	38.1	50.8	58.7	60.3	60.3	57.1	54.0	49.2	39.7	31.7	25.4
4L5N	54.2	37.5	58.3	66.7	66.7	70.8	70.8	70.8	66.7	66.7	66.7	62.5
4LGR	70.6	41.2	58.8	64.7	64.7	64.7	70.6	64.7	64.7	58.8	47.1	41.2
4M1G	30.6	36.5	47.1	54.1	55.3	54.1	51.8	47.1	42.4	37.6	31.8	27.1
4MA7	63.3	36.7	50.0	60.0	60.0	63.3	63.3	63.3	60.0	53.3	50.0	43.3
4N6R	27.9	39.5	51.2	58.1	58.1	58.1	55.8	53.5	46.5	41.9	34.9	30.2
4N90	40.0	40.0	52.0	56.0	60.0	60.0	60.0	56.0	56.0	48.0	44.0	40.0
4NBX	53.6	42.9	53.6	60.7	64.3	64.3	60.7	60.7	57.1	50.0	46.4	42.9
4NBY	65.2	39.1	56.5	65.2	69.6	69.6	69.6	69.6	65.2	56.5	52.2	47.8
4NRH	57.1	42.9	71.4	85.7	85.7	100.0	100.0	100.0	100.0	100.0	100.0	100.0
4P78	46.2	46.2	53.8	69.2	69.2	69.2	69.2	69.2	69.2	61.5	53.8	46.2
4PJ2	44.4	40.7	51.9	59.3	63.0	59.3	59.3	59.3	55.6	44.4	33.3	25.9
4QAF	60.9	34.8	47.8	54.3	56.5	56.5	54.3	52.2	47.8	41.3	32.6	26.1
4QLP	23.5	34.5	42.9	47.9	48.7	47.1	43.7	38.7	33.6	27.7	21.8	17.6
4RGO	27.5	37.3	49.0	54.9	54.9	54.9	54.9	52.9	47.1	43.1	37.3	31.4
4TSB	46.3	38.8	50.7	58.2	59.7	58.2	56.7	53.7	46.3	38.8	32.8	25.4
4TXV	44.0	40.0	52.0	56.0	56.0	56.0	56.0	52.0	48.0	44.0	40.0	36.0
4UHP	33.3	38.9	55.6	61.1	61.1	61.1	55.6	55.6	55.6	50.0	50.0	44.4
4WEM	50.0	41.2	52.9	61.8	61.8	61.8	61.8	55.9	52.9	47.1	44.1	41.2
4WKZ	37.5	41.7	58.3	62.5	66.7	66.7	62.5	58.3	54.2	45.8	41.7	37.5
4YDY	60.0	40.0	60.0	65.0	70.0	70.0	70.0	65.0	65.0	55.0	50.0	45.0
5A6W	50.0	42.9	57.1	64.3	64.3	67.9	67.9	64.3	64.3	60.7	57.1	53.6
5BV7	29.6	40.7	55.6	59.3	63.0	63.0	63.0	55.6	48.1	40.7	29.6	18.5
5BVP	50.0	40.0	53.3	63.3	63.3	63.3	63.3	63.3	60.0	56.7	50.0	46.7
5C50	51.7	37.9	55.2	62.1	62.1	62.1	58.6	55.2	55.2	51.7	48.3	44.8
5CEC	51.6	40.3	54.8	62.9	62.9	62.9	61.3	56.5	50.0	43.5	38.7	32.3

5CZX	42.1	36.8	49.1	56.1	57.9	57.9	56.1	54.4	49.1	43.9	38.6	29.8
5D2M	33.3	33.3	41.7	45.8	45.8	45.8	45.8	41.7	41.7	37.5	33.3	29.2
5D93	38.7	37.1	48.4	54.8	54.8	54.8	53.2	51.6	50.0	48.4	45.2	41.9
5DC4	54.2	41.7	58.3	66.7	66.7	66.7	66.7	66.7	58.3	54.2	50.0	37.5
5DJT	43.8	37.5	50.0	56.3	56.3	56.3	53.1	46.9	40.6	34.4	31.3	28.1
5DOI	30.8	30.8	38.5	46.2	46.2	46.2	46.2	46.2	38.5	38.5	30.8	23.1
5EE4	36.0	36.0	48.0	56.0	56.0	56.0	56.0	52.0	48.0	44.0	40.0	32.0
5ELU	75.0	41.7	58.3	66.7	75.0	75.0	75.0	75.0	66.7	66.7	58.3	58.3
5F1K	71.4	42.9	57.1	66.7	66.7	71.4	71.4	66.7	66.7	61.9	57.1	52.4
5F72	36.0	36.0	48.0	52.0	56.0	52.0	52.0	48.0	44.0	40.0	32.0	28.0
5GGS	42.1	36.8	45.6	52.6	52.6	52.6	49.1	47.4	43.9	40.4	36.8	29.8
5H5Z	47.1	38.2	52.9	61.8	61.8	61.8	61.8	58.8	55.9	52.9	47.1	41.2
5IWB	34.4	40.6	50.0	56.3	59.4	59.4	59.4	56.3	56.3	50.0	40.6	34.4
5J3T	47.1	41.2	58.8	64.7	70.6	70.6	70.6	64.7	64.7	58.8	47.1	41.2
5JDS	47.4	36.8	52.6	57.9	57.9	57.9	57.9	57.9	52.6	47.4	42.1	36.8
5KVF	32.7	36.7	48.0	55.1	56.1	55.1	52.0	48.0	41.8	35.7	28.6	22.4
5KW9	60.0	40.0	55.0	65.0	65.0	65.0	65.0	65.0	60.0	60.0	55.0	55.0
5L21	34.1	39.0	51.2	56.1	58.5	58.5	56.1	56.1	53.7	51.2	48.8	46.3
5LB7	58.3	37.5	50.0	58.3	62.5	62.5	62.5	58.3	54.2	45.8	29.2	16.7
5M20	46.7	43.3	50.0	56.7	56.7	56.7	53.3	46.7	40.0	33.3	26.7	23.3
5002	44.4	38.9	50.0	61.1	61.1	61.1	55.6	55.6	50.0	50.0	50.0	44.4
5V1Y	40.0	36.0	48.0	52.0	52.0	52.0	52.0	48.0	48.0	44.0	40.0	36.0
5VAG	44.4	38.9	55.6	61.1	61.1	66.7	66.7	61.1	61.1	55.6	50.0	38.9
5WCA	31.5	32.4	40.5	46.8	46.8	45.9	45.0	40.5	36.0	32.4	27.9	25.2
5WVO	30.8	38.5	46.2	53.8	53.8	53.8	53.8	53.8	50.0	50.0	46.2	42.3
6APP	38.9	38.9	50.0	55.6	58.3	58.3	55.6	52.8	50.0	47.2	44.4	41.7
6AZZ	42.9	42.9	57.1	66.7	66.7	71.4	71.4	66.7	61.9	57.1	52.4	52.4