## i Imię i nazwisko

#### Dominik Gront

# 2 Posiadane dyplomy, stopnie naukowe/ artystyczne – z podaniem nazwy, miejsca i roku ich uzyskania oraz tytułu rozprawy doktorskiej

23 czerwca 2001	Magister chemii, Uniwersytet Warszawski, Wydział Che- mii, Warszawa <i>"Klasyczne i nowe techniki Monte Carlo w termodynamice białek"</i>
16 kwietnia 2006	Doktor nauk chemicznych w zakresie chemii teore- tycznej, Uniwersytet Warszawski, Wydział Chemii, War- szawa "Opracowanie algorytmu do modelowania struktur białkowych na podstawie baz danych sekwencji i struktur"

# 3 INFORMACJE O DOTYCHCZASOWYM ZATRUDNIENIU W JEDNOSTKACH NAUKOWYCH / ARTYSTYCZNYCH

październik 2006 - wrzesień 2007	Uniwersytet Warszawski, Wydział Chemii, Warszawa
październik 2007 - wrzesień 2008	Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA
od października 2008	Uniwersytet Warszawski, Wydział Chemii, Warszawa

# 4 WSKAZANIE OSIĄGNIĘCIA WYNIKAJĄCEGO Z ART. 16 UST. 2 USTAWY Z DNIA 14 MARCA 2003 R. O STOPNIACH NAUKOWYCH I TYTULE NAUKOWYM ORAZ O STOPNIACH I TYTULE W ZAKRESIE SZTUKI (DZ. U. NR 65, POZ. 595 ZE ZM.)

4.1 tytuł osiągnięcia naukowego

"Opracowanie nowych algorytmów modelowania białek i ich implementacja w pakiecie BioShell"

- 4.2 publikacje wchodzące w skład osiągnięcia naukowego
- H1. D. Gront\*, S. Kmiecik, A. Kolinski, "Backbone building from quadrilaterals: A fast and accurate algorithm for protein backbone reconstruction from alpha carbon coordinates.", J. Comput. Chem. 2007; 28 1593-1597
- H2. **D. Gront**\*, A. Kolinski, *"Efficient scheme for optimization of parallel tempering Monte Carlo method.*", J. Phys.: Condens. Matter, 2007; 19 036225
- H3. A. Sikorski\* and **D. Gront**, *"Thermodynamic properties of polypeptide chains. Parallel* tempering Monte Carlo simulations", Acta Physica Polonica B, 2007 38 1899-1908
- H4. **D. Gront**\*, A. Kolinski, *"T-Pile a package for thermodynamic calculations of biomolecules.*", Bioinformatics; 2007; 23: 1840 - 1842.
- H5. **D. Gront**<sup>\*</sup> and A. Kolinski, *"Comparative modeling without implicit sequence alignments*", Bioinformatics, 2007; 23 2522
- H6. **D. Gront**\*, A. Kolinski, *"Utility library for structural bioinformatics"*, Bioinformatics, 2008; 24 584
- H7. **D. Gront**\*, A. Kolinski, "A fast and accurate methods for predicting short-range constraints in protein models", J. Comput. Aided Mol. Des, 2008 DOI 10.1007/s10822-008-9213-8
- H8. P. Gniewek, A. Kolinski, D. Gront, "Optimization of profile-to-profile alignment parameters for one-dimensional threading", J.Comput. Biol., 2012 19:879-886, doi: 10.1089/cmb.2011.0307
- H9. **D. Gront**\*, P. Wojciechowski, M. Blaszczyk, A. Kolinski, *"Bioshell Threader: protein homology detection based on sequence profiles and secondary structure profiles*", Nucl. Acid Res., 2012, doi: 10.1093/nar/gks555
- H10. **D. Gront**\*, S. Kmiecik, M. Blaszczyk, D. Ekonomiuk, and A. Kolinski, *"Optimiza-tion of protein models*", WIREs Comput Mol Sci, 2012; 2 479-493 doi: 10.1002/wcms.1090
- H11. P. Gniewek, A. Kolinski, A. Kloczkowski, **D. Gront**\*, *"BioShell-Threading: versatile Monte Carlo package for protein threading*" BMC Bioinformatics 2014 15:22
- H12. L. Wieteska<sup>\*</sup>, M. Ionov, J. Szemraj, A. Kolinski, C. Feller, **D. Gront**<sup>\*</sup> *"Improving thermal stability of thermophilic L-threonine aldolase from Thermotoga maritima*" J. of Biotechnology, 2015 199:69-76, doi: 10.1016/j.jbiotec.2015.02.013

4.3 omówienie celu naukowego/artystycznego ww. pracy/prac i osiągniętych wyników wraz z omówieniem ich ewentualnego wykorzystania.

## WPROWADZENIE

Postęp w nauce i rozwój technologii są nierozerwalnie ze sobą związane. Nauki podstawowe tworzą nowe narzędzia badawcze, które z kolei otwierają przed człowiekiem nieznane obszary badań. Pierwsze maszyny cyfrowe, powstałe w latach pięćdziesiątych minionego stulecia, przyczyniły się do intensywnego rozwoju istniejących dziedzin nauki. Zainicjowały też powstanie kilku nowych, przede wszystkim informatyki. U jej podwalin legły zarówno prace teoretyczne, jak i praktyczne rozwiązania: architektura von Neumanna<sup>[1]</sup> czy pierwsze języki programowania<sup>[2]</sup>. W latach pięćdziesiątych XX wieku opublikowane zostały dwie podstawowe metody modelowania molekularnego: dynamika molekularna<sup>[3]</sup> (na razie w ujęciu dyskretnym) oraz schemat Metropolisa<sup>[4]</sup>.

W tych samych latach postępował rozwój biologii molekularnej. Rozszyfrowano kod genetyczny<sup>[5]</sup>, poznano strukturę DNA<sup>[6]</sup> oraz mioglobiny<sup>[7]</sup>. Opisano także zasady rządzące strukturą białka<sup>[8–10]</sup>. Zautomatyzowano też metodę sekwencjonowania białek, co zaowocowało opublikowaniem - w formie książki - przez Margaret Dayhoff pierwszej sekwencyjnej bazy danych<sup>[11]</sup>.

Badania te możliwe były również dzięki wykorzystaniu pierwszych komputerów. Zastosowano je między innymi do rozwiązywania rentgenograficznych struktur białek a także do składania fragmentarycznych wyników sekwencjonowania<sup>[12]</sup>. Sekwencji białek znano już coraz więcej; w literaturze zaczęto dyskutować więc ewolucję genów i białek. W roku 1967 opublikowano pierwsze drzewo filogenetyczne<sup>[13]</sup>. Wreszcie w 1970r. Needleman i Wunsch opublikowali algorytm globalnego uliniowienia sekwencji<sup>[14]</sup>. W ten sposób zrodziła się nowa nauka - *bioinformatyka* - choć na samą jej nazwę trzeba było jeszcze poczekać prawie dziesięć lat<sup>[15]</sup>.

<sup>[11]</sup> Dayhoff, M. O. (Silver Spring, Md., 1965).

<sup>&</sup>lt;sup>[1]</sup> von Neumann, J. Tech. rep. (Philadelphia, PA, USA, 1945), 1-43.

<sup>&</sup>lt;sup>[2]</sup> Backus, J. W. et al. in Papers Presented at the February 26-28, 1957, Western Joint Computer Conference: Techniques for Reliability (Los Angeles, California, 1957), 188–198.

<sup>&</sup>lt;sup>[3]</sup> Alder, B. J. & Wainwright, T. E. The Journal of Chemical Physics 27, 1208–1209 (1957).

<sup>[4]</sup> Metropolis, N. et al. The Journal of Chemical Physics 21, 1087–1092 (1953).

<sup>[5]</sup> Gamow, G. et al. Advances in biological and medical physics 4, 23–68 (1956).

<sup>&</sup>lt;sup>[6]</sup> Watson, J. D. & Crick, F. H. Nature **171**, 737–738 (1953).

<sup>&</sup>lt;sup>[7]</sup> Kendrew, J. C. et al. Nature **181**, 662–666 (1958).

<sup>[8]</sup> Pauling, L. & Corey, R. B. Proceedings of the National Academy of Sciences 37, 251–256 (1951).

<sup>&</sup>lt;sup>[9]</sup> Pauling, L. et al. Proceedings of the National Academy of Sciences **37**, 205–211 (1951).

<sup>&</sup>lt;sup>[10]</sup> Ramachandran, G. N. et al. Journal of molecular biology 7, 95–99 (1963).

<sup>&</sup>lt;sup>[12]</sup> Dayhoff, M. O. & Ledley, R. S. in *Proceedings of the December 4-6, 1962, Fall Joint Computer Conference* (Philadelphia, Pennsylvania, 1962), 262–274.

<sup>&</sup>lt;sup>[13]</sup> Fitch, W. M. & Margoliash, E. Science (New York, N.Y.) 155, 279–284 (1967).

<sup>&</sup>lt;sup>[14]</sup> Needleman, S. B. & Wunsch, C. D. *Journal of molecular biology* **48**, 443–453 (1970).

<sup>&</sup>lt;sup>[15]</sup> Hogeweg, P. PLoS Comput Biol 7, e1002021+ (2011).

Kierunki rozwoju tej dziedziny w kolejnych dekadach związane były z gwałtownie powiększającymi się zasobami danych: sekwencjami i strukturami biomolekuł. Zadaniem bioinformatyki stało się te dane gromadzić i analizować. By temu sprostać, rozpoczęto tworzenie odpowiedniego oprogramowania. Równocześnie zaczęło powstawać oprogramowanie do modelowania struktury i dynamiki biomolekuł. Początkowo oba te obszary były wyraźnie rozgraniczone. Bioinformatyka bazowała na danych, modelowanie molekularne zaś - na zasadach fizyki. Okazało się jednak, że wnioskowanie oparte o ewolucję genów może być równie ważne co modele wyprowadzone z pierwszych zasad.

Aktualnie istnieje wiele inicjatyw poświęconych tworzeniu oprogramowania zarówno z dziedziny bioinformatyki jak i modelowania molekularnego. Spośród narzędzi do modelowania należy wspomnieć pakiety do dynamiki molekularej, programy Modeller<sup>[16]</sup>, ICM<sup>[17]</sup>, UNRES<sup>[18]</sup>, Rosetta<sup>[19]</sup>, czy rodzinę modeli siatkowych opracowanych w grupie prof. Kolińskiego (np. SICHO<sup>[20]</sup> oraz CABS<sup>[21]</sup>). Z kolei jako przykłady typowych pakietów bioinformatycznych wymienić należy BioPerl<sup>[22]</sup>, BioPython<sup>[23]</sup>, BioJava<sup>[24]</sup> i BioRuby<sup>[25]</sup>. W roku 2004 pakiety te<sup>1</sup>, udostępniały jedynie bardzo skromny zestaw funkcji operujących na strukturach białek. Funkcjonalność taka była jednakże niezbędna piszącemu te słowa do sprawnego realizowania wykonywanej przezeń pracy doktorskiej. Było to impulsem, który zainicjował powstanie pakietu BioShell.

<u>Celem naukowym przedstawionego cyklu prac</u> było stworzenie spójnego i kompletnego pakietu oprogramowania, wspomagającego rozwiązywanie różnorodnych problemów z zakresu bioinformatyki strukturalnej oraz modelowania struktur biomolekuł. Praktycznie cały kod źródłowy napisany został przez habilitanta a powstałe narzędzia obliczeniowe wykorzystywane są w co najmniej kilkunastu laboratoriach na świecie. W pracach tych przedstawiono konstrukcję kolejnych wersji pakietu oraz opisano najważniejsze zaimplementowane w nim algorytmy. Do cyklu zostały włączone również prace prezentujące przykładowe zastosowania pakietu.

<sup>&</sup>lt;sup>1</sup> jedynie BioPerl i BioPython; BioJava i BioRuby jeszcze wtedy nie istniały

<sup>&</sup>lt;sup>[16]</sup> Šali, A. & Blundell, T. L. *Journal of Molecular Biology* 234, 779-815 (1993).

<sup>&</sup>lt;sup>[17]</sup> Abagyan, R. et al. J. Comput. Chem. 15, 488-506 (1994).

<sup>&</sup>lt;sup>[18]</sup> Liwo, A. et al. The Journal of Chemical Physics **115**, 2323–2347 (2001).

<sup>&</sup>lt;sup>[19]</sup> Rohl, C. A. et al. in Numerical Computer Methods, Part D 66-93 (Department of Biochemistry and Ho-

ward Hughes Medical Institute, University of Washington, Seattle, Washington 98195, USA., 2004).

<sup>&</sup>lt;sup>[20]</sup> Kolinski, A. & Skolnick, J. Proteins **32**, 475–494 (1998).

<sup>&</sup>lt;sup>[21]</sup> Kolinski, A. Acta biochimica Polonica **51**, 349–371 (2004).

<sup>&</sup>lt;sup>[22]</sup> Stajich, J. E. *et al. Genome research* **12**, 1611–1618 (2002).

<sup>&</sup>lt;sup>[23]</sup> Chapman, B. & Chang, J. SIGBIO Newsl. 20, 15–19 (2000).

<sup>&</sup>lt;sup>[24]</sup> Holland, R. C. G. et al. Bioinformatics 24, 2096–2097 (2008).

<sup>&</sup>lt;sup>[25]</sup> Goto, N. et al. Bioinformatics 26, 2617-2619 (2010).

# Oprogramowanie BioShell

#### Wersja 1.x - programy uruchamiane z linii poleceń systemu UNIX

Konstrukcja pierwszej wersji pakietu, opublikowanej w 2006 roku<sup>[26]</sup>, w całości opierała się na idei funkcjonowania systemu UNIX. Składał się on więc z kilku programów których działanie kontrolowane było odpowiednimi opcjami podawanymi z linii poleceń. W początkowej wersji, BioShell ułatwiał prowadzenie symulacji dynamiki białek w modelu zredukowanym CABS<sup>[21]</sup>. Wykorzystywany był do przygotowywania plików wejściowych i analizowania wynikowych trajektorii. Inne moduły pakietu posłużyły do wyprowadzania potencjałów średniej siły na podstawie statystyk zebranych ze znanych struktur białkowych. W roku 2007 na pakiet składały się następujące programy:

- strc (**str**ucture **c**onverter) dokonujący konwersji pomiędzy formatami plików zapisujących struktury biomolekuł
- str\_calc (structure calculator) wykonujący obliczenia na strukturach białek: mapy kontaktów, kąty łańcucha głównego  $\Phi$ ,  $\Psi$ ,  $\omega$  oraz łańcuchach bocznych  $\chi$ , itp.
- rms\_calc (rms calculator) obliczający optymalne nałożenie jednej struktury białka na drugą.
  - clust (clustering) do analizy skupień
  - alignc (**align**ment **c**onverter) dokonujący konwersji pomiędzy formatami uliniowień sekwencyjnych
- proline (**pr**ofile **aligne**r) służący do znajdowania optymalnego uliniowienia dwóch sekwencji lub profili sekwencyjnych.

Duży nacisk położono na integrację pakietu ze standardowymi poleceniami systemu UNIX, takimi jak grep, sed czy awk.

#### Wersja 2.x - biblioteka modułów dla języków skryptowych

Oprogramowanie opisane powyżej doskonale spełniało swoje zadanie, ale jego rozbudowa o nowe funkcje nastręczała kilku trudności. Najpoważniejszą z nich było jednoznaczne a zarazem elastyczne zdefiniowanie kolejności wykonywania zadanych operacji. Dlatego po uzyskaniu stopnia doktora, autor rozpoczął prace nad kolejną wersją pakietu. Diametralnej zmianie uległa filozofia, na której oparto konstrukcję oprogramowania. W nowej wersji BioShell stał się przede wszystkim biblioteką funkcji, wywoływanych ze skryptów napisanych w języku Python<sup>H6.</sup>. Rozwiązało to całkowicie problem sterowania obliczeniami oraz znacznie poszerzyło wachlarz funkcji pakietu udostępnionych użytkownikowi. W nowej wersji pakietu odtworzono programy z wersji poprzedniej. Dodano też dwa nowe: PsiBlastSearch i PsiBlastAnalyse<sup>2</sup>, służące do analizowania przestrzeni

<sup>&</sup>lt;sup>2</sup> Wersja druga pakietu napisana została w języku Java. Nazwy wszystkich programów pakietu zmieniono, wprowadzając wielkie litery zgodnie z ogólnie przyjętymi w tym środowisku konwencją.

<sup>&</sup>lt;sup>[26]</sup> Gront, D. & Kolinski, A. *Bioinformatics* 22, 621-622 (2006).



sekwencji białkowych w pobliżu zadanej sekwencji-celu.

Rysunek I: Hierarchiczna struktura biblioteki BioShell (określonej tu jbcl -Java BioComputing Library). Poszczególne moduły pogrupowane zostały w pakiety odpowiadające ich funkcjonalności. Dla przykładu algorytmy operujące na grafach zawarto w jbcl.algorithms.graphs. Procedury do wyznaczania ulinowień - w jbcl.calc.alignment, a same funkcje oceniające uliniowienia - w jbcl.calc.alignment.scoring.

Oprogramowanie dostępne jest na stronie bioshell.pl, wraz z obszerną dokumentacją, przykładowymi skryptami, danymi testowymi itp.

#### #!/usr/bin/env jython

import sys # Here we import two BioShell modules: PDB (to read PDB files) ... from jbcl.data.formats import PDB # ... and Neighborhood to look for spatial neighbours. from jbcl.calc.structural import Neighborhood inputFile = sys.argv[1] # PDB file name is the parameter of this script reader = PDB(inputFile) allAtoms = reader.getStructure().getAtomsArray()

n = Neighborhood(allAtoms)
cuResidues = protein.findResidues("\_CU\_")

for cuResidue in cuResidues :
 cuAtom = cuResidue.getAtomsArray()[0]
 nn = n.findNeighborsArray(cuAtom,4.0)
 residueSet = Set()
 for atom in nn : residueSet.add( atom.getOwner() )

for residue in residueSet :
 for line in PDB.createPdbLines( residue ) : print line



Rysunek 2: **Przykładowy skrypt** wykorzystujący bibliotekę BioShell *(po lewej)*. Skrypt ten ładuje plik PDB i wyszukuje wszystkie atomy miedzi (rozpoznawane po nazwie atomu). Dla każdego z tych atomów wyszukuje jego otoczenie w przestrzeni - reszty aminokwasowe odległe o nie więcej niż 4Å- a wyniki nagrywane są w formacie PDB. *Powyżej*: przykładowy fragment struktury wycięty z depozytu 2AZA (azuryna) za pomocą w/w skryptu.

#### ZASTOSOWANIA

#### Modelowanie struktur białek

Od wielu lat głównym zainteresowaniem naukowym autora pozostaje modelowanie struktur białek. Większość modułów pakietu powstała w tym właśnie celu. Procedury te bezpośrednio wykonują potrzebne obliczenia, bądź automatyzują pracę zewnętrznych programów, takich jak PsiBlast, Modeller<sup>[16]</sup> czy Rosetta<sup>[19]</sup>. Pakiet BioShell wykorzystywano kilkukrotnie podczas odbywających się co dwa lata eksperymentów CASP<sup>3</sup> w latach: 2004 (CASP6), 2006 (CASP7), 2010 (CASP9) i 2014 (CASP11). Eksperyment ten ma charakter konkursu, w którym uczestniczą grupy teoretyczne, starając się jak najdokładniej przewidzieć struktury białek wyznaczone eksperymentalnie. Faktyczne struktury są ujawnione dopiero po zakończeniu konkursu. Wykorzystywany przez autora protokół modelowania zmieniał się istotnie przez te lata wraz ze zdobywanym doświadczeniem i implementacją nowych algorytmów. Na Ryc. 3 przedstawiono schemat postępowania przyjęty podczas CASP11<sup>[27]</sup>. Niektóre z wykorzystanych w nim procedur obliczeniowych opisano szerzej w dalszej części autoreferatu, pokazując jednocześnie rolę oprogramowania BioShell. W końcowym rankingu<sup>4</sup> kategorii modelowania w oparciu o szablon<sup>5</sup>, uwzględniającym 81 przewidzianych domen grupa BioShell-server została sklasyfikowana na 41 pozycji. Startująca w kategorii **FM** (*Free Modelling*) grupa BioShell zajęła 40 pozycję. Ogółem z konkursie CASP11 wzięło udział 123 grupy eksperckie i 84 serwery



Rysunek 3: Schemat modelowania struk**tur białek** Diagram (zaadaptowany z pracy<sup>[27]</sup>) obrazuje schemat modelowania struktur białek, jaki wykorzystano w trakcie eksperymentu CASP11. Algorytm startuje z sekwencji białkacelu. Pierwszym etapem jest przeglądanie baz danych w poszukiwaniu białek o podobnych sekwencjach (najprawdopodobniej homologicznych). Będą one potencjalnymi szablonami do modelowania. Na podstawie stworzonego uliniowienia wielu sekwencji przewidywane są również pewne cechy strukturalne, np. struktura drugorzędowa, czy wyeksponowanie grup bocznych badanego białka do rozpuszczalnika. Informacja ta wykorzystywana jest do obliczenia uliniowienia pomiędzy sekwencją białka modelowanego a szablonami. Uliniowienia te są następnie udokładniane (przewlekanie 3D), a na ich podstawie budowane są modele strukturalne. Ostatni etap to analiza i selekcja uzyskanych modeli.

W ostatnich latach pakiet wykorzystywany jest również w zadaniach związanych z racjonalną inżynierią białek. Oba te cele są dość podobne i wykorzystują te same metody obliczeniowe. W pierwszym jednak znana jest sekwencja białka a celem modelowania jest jego struktura. W drugim zaś zakłada się konkretną strukturę trójwymiarową białka a poszukuje sekwencji.

<sup>&</sup>lt;sup>3</sup> Critical Assessment of Protein Structure Prediction Methods; http://predictioncenter.org/

<sup>&</sup>lt;sup>4</sup> http://predictioncenter.org/casp11/zscores\_final.cgi

<sup>&</sup>lt;sup>5</sup> **TBM**, Template Based Modelling

<sup>&</sup>lt;sup>[27]</sup> Strumillo, M. et al. in. **18** (2014), 379–384.

Modelowanie molekularne to nie jedyna rola pakietu. Spora jego część poświęcona jest bowiem analizie sekwencji i struktur białek. Pakiet udostępnia też wiele standardowych procedur numerycznych i statystycznych. Przykładowe zastosowania również podsumowano poniżej.

## Wyszukiwanie sekwencji

Standardowo do przeszukiwania sekwencyjnych baz danych wykorzystywany jest program PSIBLAST. Stosowany był on także w pracach podsumowanych w niniejszym Autoreferacie. Program ten doskonale sprawdza się w przypadku, gdy odległość ewolucyjna pomiędzy sekwencjami nie jest duża. W innych przypadkach najlepiej sprawdza się wielokrotne, często iteracyjne wykorzystanie tego programu<sup>[28]</sup>. W tym właśnie celu powstały programy PsiBlastSearch i PsiBlastAnalyse napisane w niecałe dwa tygodnie latem 2010 r, w trakcie trwającego właśnie konkursu CASP9. W ten sposób autor niniejszego Autoreferatu w pełni zautomatyzował wykonywaną dotychczas manualnie w laboratorium prof. Bakera procedurę, stanowiącą pierwszy etap modelowania porównawczego. Program PsiBlastSearch automatyzuje pracę narzędzia PSIBLAST, uruchamiając go z różnymi wartościami parametrów startowych. Program PsiBlastAnalyse przetwarza i analizuje uzyskane wyniki. Analiza ta obejmuje filtrowanie znalezionych sekwencji wg wielu kryteriów oraz analizę skupień. Ostatecznym wynikiem jest różnorodny<sup>6</sup> zbiór sekwencji homologicznych do sekwencji modelowanego białka, wykorzystany następnie do podziału sekwencji celu na (potencjalne) domeny oraz przy konstrukcji profilu sekwencyjnego. Procedura ta była intensywnie wykorzystywana w trakcie eksperymentu CASP9, została też zastosowana przy projektowaniu mutantów białka aldolazy treoninowej (praca H12.).

#### Predykcje w oparciu o sekwencje

Profil sekwencyjny obliczony w opisanym powyżej etapie, wykorzystywany jest przez szereg narzędzi, których celem jest predykcja pewnych cech strukturalnych nieznanego białka, np jego struktury drugorzędowej oraz wyeksponowania poszczególnych reszt do rozpuszczalnika. BioShell automatyzuje pracę programów: PsiPred, Porter, SAM, Juffo i SpineX. Uzyskane wyniki służą następnie jako dane wejściowe do tworzenia uliniowienia oraz na etapie budowy modelu.

#### Tworzenie i optymalizacja uliniowień

Niezbędnym elementem modelowania porównawczego białka-celu jest znalezienie szablonu, czyli wystarczająco podobnego białka (w założeniu: homologa bądź analoga strukturalnego), którego struktura została już wyznaczona eksperymentalnie. Bardzo istotne

<sup>&</sup>lt;sup>6</sup> ang. *non-redundant* 

<sup>&</sup>lt;sup>[28]</sup> Margelevicius, M. & Venclovas, C. BMC Bioinformatics 6, 185+ (2005).

jest również zbudowanie poprawnego uliniowienia tych dwóch białek. W piśmiennictwie znaleźć można bardzo wiele metod rozwiązywania tego problemu. Ogólnie rzecz ujmując, można je podzielić na cztery grupy: *(i)* Uliniowienie sekwencyjne pary białek, *(ii)* ulinowienie sekwencyjne całej rodziny, *(iii)* uliniowienie dwóch profili sekwencyjnych oraz *(iv)* uliniowienie sekwencji celu ze strukturą szablonu.

Metoda (i) jest mało dokładna i sprawdza się tylko w przypadku, gdy białka są bardzo blisko spokrewnione. Aby zastosować podejście (ii), trzeba jednoznacznie wskazać sekwencje białek, należące do badanej rodziny. W rękach eksperta metoda ta daje doskonałe wyniki; najczęściej wymaga jednak manualnej ingerencji w tworzone uliniowienie<sup>[29]</sup>. Ponieważ jednym z celów postawionych przed oprogramowaniem BioShell była jak najpełniejsza automatyzacja procesu modelowania biomolekuł, w pakiecie zaimplementowano łatwe do zautomatyzowania metody (iii) oraz (iv), zwane przewlekaniem białek. Pierwszą z nich - przewlekanie jednowymiarowe - opublikowano w roku 1987<sup>[30]</sup>, drugą pięć lat później<sup>[31]</sup>. Wtedy też pojawiła się nazwa "przewlekanie"<sup>7</sup>. Dla porównania, program Psi-Blast<sup>[32]</sup> opublikowano w roku 1997. Wspomniane dwa warianty przewlekania różnią się diametralnie. O ile przypadek jednowymiarowy to "zwykłe" uliniowienie dwóch profili sekwencyjnych obliczane metodą dynamicznego programowania, to wariant trójwymiarowy jest problemem NP-zupełnym. W literaturze zaproponowano kilka przybliżonych metod jego rozwiązywania. Wszystkie one wymagają jednak ogromnych zasobów obliczeniowych. Z drugiej strony, wspomniany już program PsiBlast umożliwił szybkie wybieranie z baz danych sekwencji podobnych do białka-celu i stał się wygodnym narzędziem do tworzenia profili sekwencyjnych. W konsekwencji istotnie uprościł procedury przewlekania jednowymiarowego. Czynniki te spowodowały, że algorytmy uliniowienia profili wyparły prawie zupełnie metody przewlekania trójwymiarowego. Jeszcze do niedawna jedynym publicznie dostępnym programem 3D był Raptor<sup>[33]</sup>. Metoda ta, wysoko oceniana w kolejnych konkursach przewidywania struktury białek CASP, wykorzystuje algorytm programowania liniowego.

Metoda BioShell Threading 3D<sup>HII</sup> oparta jest na zupełnie innej zasadzie. Problem optymalizacji uliniowienia potraktowany został jak zagadnienie z dziedziny modelowania molekularnego. Ocenę uliniowienia (*ang. score*) zastąpiła energia, wykorzystująca zarówno człony sekwencyjne jak i strukturalne. Algorytmem uliniowienia stało się probkówanie przestrzeni stanów układu. Stany te - czyli uliniowienia - zapisane są jako lista ciągłych (niezawierających przerw) bloków. Ruchy Monte Carlo polegają na dzieleniu, łączeniu, przesuwaniu, skracaniu oraz wydłużaniu bloków. Algorytm ten jest bardzo

<sup>&</sup>lt;sup>7</sup> ang. *threading* 

<sup>&</sup>lt;sup>[29]</sup> Venclovas, C. & Margelevicius, M. Proteins 77 Suppl 9, 81-88 (2009).

<sup>&</sup>lt;sup>[30]</sup> Gribskov, M. et al. Proceedings of the National Academy of Sciences of the United States of America **84**, 4355–4358 (1987).

<sup>&</sup>lt;sup>[31]</sup> Jones, D. T. et al. Nature **358**, 86–89 (1992).

<sup>&</sup>lt;sup>[32]</sup> Altschul, S. F. et al. Nucleic Acids Research 25, 3389–3402 (1997).

<sup>&</sup>lt;sup>[33]</sup> Xu, J. et al. J Bioinform Comput Biol 1, 95–117 (2003).

ogólny. Jako dane wejściowe można użyć dowolne połączenie sekwencji, profilu sekwencyjnego albo struktury. Dla przykładu, uruchomienie algorytmu dla pary profili jest równoważne przewlekaniu ID a dla pary struktur - zagadnieniu uliniowienia struktur. To ostatnie zagadnienie posłużyło do przetestowania opisywanej metody. W pracy **HII.** pokazano, że BioShell Threading 3D znajduje lepsze uliniowienia struktur niż tm-align - jedno z najlepszych programów stworzonych do tego celu.

Rozwiązania zaimplementowane w programie BioShell Threading 3D oczywiście nie rozwiązują ściśle problemu NP, a próbkowanie przestrzeni uliniowień wymaga znacznych zasobów. Dlatego stosowana metodyka zakłada wykorzystanie dwóch programów: przewlekania jednowymiarowego oraz przewlekania trójwymiarowego. Pierwszy z nich wykorzystuje programowanie dynamiczne do uliniowienia profili sekwencyjnych obliczonych dla białek: szablonu i celu. Profile są uzupełnione informacją o strukturze drugorzędowej, co istotnie poprawia czułość wyszukiwania<sup>H8</sup>. Metodę tą zoptymalizowano tak, aby z jak największą wiarygodnością wskazywała białka (potencjalne szablony) należące do tej samej rodziny SCOP, co białko-cel. Obliczenia przewlekania trójwymiarowego prowadzone są tylko dla potencjalnych szablonów wyłonionych w poprzednim etapie i mają na celu uzyskanie jak najlepszego uliniowienia. Bardzo istotny jest również fakt, że próbkowanie przestrzeni stanów dostarcza także wysoko ocenionych uliniowień sub-optymalnych. Na podstawie każdego z nich budowany jest strukturalny model białka-celu. Procedura ta wykorzystana została przez dwa zespoły badawcze, biorące udział w konkursie CASPII: (grupy BioShell-server i BioShell-human) oraz zespół dra Chena Keasara<sup>8</sup> (grupa keasar).

Warto zaznaczyć, że uliniowienia suboptymalne mogą być generowane także innymi metodami. Najczęściej uzyskuje się je odpowiednio modyfikując algorytm odtwarzania ścieżki uliniowienia<sup>9</sup>. Do pakietu BioShell dołączono implementację eleganckiego algorytmu Miyazawy<sup>[34]</sup>, który zachowuje rozkład Boltzmanna generowanych uliniwień. Algorytm ten jednakże, jako wariant programowania dynamicznego, nie umożliwia wykorzystania pełnej informacji o trójwymiarowe strukturze szablonu. Dlatego też został zastąpiony przez opisany powyżej program BioShell Threading 3D.

Zupełnie inne podejście do modelowania porównawczego zaproponowano w pracy **H5**.; w podejściu tym wejściowe uliniowienie pomiędzy modelowanym białkiem a jego szablonem w zasadzie nie jest potrzebne. Przestrzenna struktura szablonu zrzutowana jest na siatkę w modelu CABS. Warto w tym miejscu zaznaczyć, że w tym modelu wszystkie atomy węgli  $\alpha$  leżą na sieci kubicznej o stałej 0.61Å. Do modelu wprowadzono dodatkowy człon energii, oceniający dopasowanie pomiędzy aktualną modelowaną konformacją a siatkowym rzutem szablonu. Nagroda (bądź kara) energetyczna przyznawana jest,

<sup>&</sup>lt;sup>8</sup> Ben-Gurion University of the Negev, Be'er Sheva, Israel

<sup>&</sup>lt;sup>9</sup> ang. *backtracking* 

<sup>&</sup>lt;sup>[34]</sup> Miyazawa, S. Protein Eng. 8, 999–1009 (1995).

gdy atomy szablonu i celu znajdują się w tych samych węzłach sieci. Metoda ta doskonale sprawdza się w przypadku modelowania małych białek, jednakże dla dużych molekuł wymaga znacznych zasobów obliczeniowych.

#### Budowa modelu

Do budowy modelu w oparciu o szablon wykorzystywany jest program Modeller, przypadkach modelowania *de novo* zaś - Rosetta. Dodatkowo w obu scenariuszach wykorzystywany jest także program CABS. Zastosowanie tego ostatniego, w którym łańcuch polipeptydowy zapisany jest w reprezentacji zredukowanej, pociąga za sobą konieczność odbudowy detali atomowych modelu. Specjalnie w tym celu powstał program BBQ<sup>HI.</sup>, który rekonstruuje szkielet główny białka. Grupy boczne odtwarzane są programem scwrl<sup>[35]</sup>.

#### Analiza skupień

Wynikiem modelowania struktur białek jest zazwyczaj bardzo wiele modeli. Kolejnym etapem jest zatem analiza skupień, której celem jest wybór reprezentatywnych struktur. W pakiecie BioShell zaimplementowano algorytm hierarchiczny<sup>[36]</sup>. Program clust<sup>10</sup> rutynowo analizuje zbiory po kilkanaście a nawet kilkadziesiąt tysięcy struktur. Narzędzie to zostało napisane na potrzeby CASP6, a wykorzystano je także w konkursach CASP7 oraz CASP11. W pracy<sup>[37]</sup> posłużył do analizowania wyników dokowania krótkiego peptydu pochodzącego z białka C3D do domeny SH3-N. Procedura hierarchiczna jest również wykorzystana w programie PsiBlastAnalyse do grupowania podobnych sekwencji białek.

#### Procedury obliczeniowe dla struktur biomolekuł oraz metody numeryczne

Pakiet BioShell oferuje bardzo szeroki wachlarz algorytmów dedykowanych do analizy struktur białek. Oblicza ich nałożenia i uliniowienia. Wyznacza też różne parametry strukturalne: kąty płaskie i torsyjne, odległości, mapy kontaktów i mapy wiązań wodorowych. Pakiet BioShell udostępnia również szeroki zasób metod numerycznych i statystycznych, służących do przetwarzania danych, np. metody interpolacji, *bootstrap*, histogramy czy też estymatory jądrowe. Funkcjonalność ta jest nieczęsto spotykana w takich pakietach, a niektóre funkcje są unikalne dla pakietu BioShell. Być może dlatego właśnie inne laboratoria wykorzystujące ten pakiet w swoich badaniach, przede wszystkim sięgają po procedury operujące na strukturach biomolekuł<sup>[38,39]</sup>.

#### Potencjały statystyczne

<sup>&</sup>lt;sup>10</sup> w oryginalne pracy opublikowano go pod nazwą HCPM - Hierarchical Clustering of Protein Models

<sup>&</sup>lt;sup>[35]</sup> Dunbrack, J. & Karplus, M. Journal of Molecular Biology 230, 543-574 (1993).

<sup>&</sup>lt;sup>[36]</sup> Gront, D. & Kolinski, A. *Bioinformatics* **21**, 3179-3180 (2005).

<sup>&</sup>lt;sup>[37]</sup> Gront, D. et al. Acta Pol Pharm **63**, 436–438 (2006).

<sup>&</sup>lt;sup>[38]</sup> Kim, H. & Kihara, D. Proteins 82, 3255-3272 (2014).

<sup>&</sup>lt;sup>[39]</sup> Chruszcz, M. et al. Journal of Biological Chemistry 287, 7388–7398 (2012).

Funkcjonalność ta została wykorzystana kilkukrotnie do wyprowadzania potencjałów statystycznych. Dla przykładu, w pracy H7 zaproponowano potencjały statystyczne opisujące lokalną (tj. obejmującą kilka następujących po sobie reszt aminokwasowych) geometrię łańcucha głównego w konkretnej rodzinie białek. Potencjały statystyczne są powszechnie wykorzystywane w modelowaniu struktur biomolekuł. W typowym ujęciu opisują one prawdopodobieństwo pojawienia się pewnej własności strukturalnej (np. kontaktu między atomami), uzależnione od typów aminokwasów. Dla przykładu, w programie CABS wykorzystywane są potencjały oceniające odległość  $R_{15}(A_2, A_3)$ , czyli między każdym *i*-tym a *i*+4-tym węglem  $\alpha$  wzdłuż łańcucha głównego białka. Oprócz wspomnianej odległości, funkcja ta zależy od typu aminokwasów na pozycjach *i*+1 oraz *i*+3. Raz wyprowadzona, może być wykorzystywana do modelowania białka o dowolnej se-kwencji.

W przypadku takich potencjałów opisujących rodzinę białek, funkcja energii zależy od pozycji w sekwencji modelowanego białka. Na tej samej zasadzie oparte są m.in. profile sekwencyjne<sup>[30]</sup> oraz biblioteki fragmentów, wykorzystywanych w modelowaniu struktur białek<sup>[40]</sup>. Podobnie jak profil czy fragmenty, potencjał taki musi być wyprowadzony oddzielnie dla każdej sekwencji modelowanego białka. Dodatkowo, bazy danych muszą zawierać informację o strukturach białek należących do tej samej rodziny, co białko badane. Uzyskany potencjał dostarcza jednak znacznie dokładniejszego opisu lokalnej konformacji łańcucha polipeptydowego.

Pakiet BioShell został także wykorzystany do analizy struktury białka AF2331<sup>[41]</sup> (depozyt w PDB: 2FD0). Białko to jest nietypowym reprezentantem klasy  $\alpha + \beta$ , bowiem jedną z dwóch jego  $\beta$ -kartek tworzą naprzemiennie fragmenty łańcuchów A oraz B. W ten sposób oddziaływania pomiędzy łańcuchami odpowiadają za jego znaczną część struktury drugorzędowej. Z ryc. 4, przedstawiającej wybrane  $\beta$ -wstęgi tego białka można odczytać, że podążając za wiązaniami wodorowymi  $\beta$ -kartki, odwiedzamy łańcuchy AABABABB, zmieniając kod łańcucha aż pięciokrotnie. Oczywistym pytaniem, które nasunęło się w trakAnaliza topologii białek splecionych<sup>11</sup>



Rysunek 4: **Spleciona**  $\beta$ -kartka białka 2FD0. Tworzy je osiem wstęg z łańcuchów A oraz B depozytu, oznaczonych odpowiednio kolorami niebieskim i czerwonym. Pomarańczowe przerywane linie ilustrują sieć wiązań wodorowych.

cie badań nad 2FD0 było, czy topologia taka pojawiła się już w znanych strukturach.

<sup>&</sup>lt;sup>11</sup> ang. *interdigitated* 

<sup>&</sup>lt;sup>[40]</sup> Gront, D. et al. PloS one **6**, e23294+ (2011).

<sup>&</sup>lt;sup>[41]</sup> Wang, S. et al. Protein Science 18, 2410–2419 (2009).

Odpowiedź na nie przyniósł krótki skrypt napisany w środowisku BioShell. Skrypt ten wczytywał plik PDB, z którego odczytywał informację o  $\beta$  wstęgach, i obliczał sieć wiązań wodorowych. Następnie budował graf, którego wierzchołkami były  $\beta$  wstęgi pokolorowane wg kodu łańcucha a krawędziami - wiązania wodorowe. Ostatecznym wynikiem była najdłuższa możliwa w tym grafie ścieżka wiodąca przez węzły o naprzemiennie różnych kolorach. Analiza przeprowadzona na wszystkich znanych ówcześnie strukturach białek pokazała, że o ile pojedynczo-splecione  $\beta$  wstęgi (układ typu ABA) są często spotykane, to już topologie ABAB (trzykrotna zmiana kodu łańcucha) obserwowano około stukrotnie. Czterokrotna zmiana łańcucha (ABABA) pojawiła się tylko w 3 depozytach. Układ ABABAB, wykryty w AF2331, pojawił się jeszcze tylko w depozycie 2HJ1 i był najdłuższym zaobserwowanym.

#### Statystyczny opis stanów w rozkładzie kanonicznym

Modelowanie biomolekuł często prowadzone jest tak, aby symulowany układ opisywany był zespołem kanonicznym. Jest to szczególnie łatwe, gdy procedura modelowania oparta jest o metodę Monte Carlo wg wspomnianego już we wstępie schematu Metropolisa. Wynikiem modelowania jest wtedy zbiór konformacji układu o energiach opisanych rozkładem Boltzmanna. W pakiecie BioShell zaimplementowano metodę ważonych bisto-gramów<sup>[42]</sup>, dzięki której możliwe jest wyznaczenie sumy statystycznej  $\mathcal{Z}(T)$  badanego układu. Danymi wejściowymi do programu MultiHist<sup>H4.</sup> pakietu BioShell są wartości energii  $\mathcal{E}$  zaobserwowane podczas symulacji w różnych temperaturach a wynikiem - suma statystyczna  $\mathcal{Z}(T)$  oraz gęstości stanów  $\Omega(\mathcal{E})$ . Program StatPhys z kolei wczytuje  $\Omega(\mathcal{E})$  oraz zmierzone observable, dla których wyznacza średnie kanoniczne w dowolnej temperaturze. Programy te wykorzystano do badania sieciowych modeli prostych polimerów<sup>H3.</sup> a także do analizowania wyników symulacji białek w zredukowanym modelu CABS<sup>H2.</sup> (metodą Monte Carlo) a także pełnoatomową dynamiką molekularną<sup>[43]</sup>.

Symulacje te prowadzone były metodą wymiany replik Monte Carlo<sup>[44]</sup> (REMC<sup>12</sup>). W metodzie tej, dzięki jednoczesnemu modelowaniu replik tego samego układu w wielu temperaturach, eksploracja przestrzeni stanów jest znacznie wydajniejsza. Wymiana replik, czyli przenoszenie kopii układu pomiędzy temperaturami, pozwala na względnie łatwe przekraczanie barier energetycznych. Kluczowy jest tu jednak odpowiedni wybór zbioru temperatur  $\{T_i\}$ , w których symulowane są poszczególne repliki. W literaturze opisano wiele rozwiązań<sup>[45,46]</sup>, żadne z nich jednak nie gwarantuje optymalnego przepływu replik pomiędzy temperaturami. W pracy **H3** zaproponowano nowatorski sposób wyboru temperatur do prowadzenia symulacji. Opiera się on na spostrzeżeniu,

<sup>&</sup>lt;sup>12</sup> Replica Exchange Monte Carlo

<sup>&</sup>lt;sup>[42]</sup> Ferrenberg, A. M. & Swendsen, R. H. Physical Review Letters 63, 1195–1198 (1989).

<sup>&</sup>lt;sup>[43]</sup> Wabik, J. et al. International Journal of Molecular Sciences 14, 9893–9905 (2013).

<sup>&</sup>lt;sup>[44]</sup> Geyer, C. J. in *Computing Science and Statistics: Proceedings of 23<sup>rd</sup> Symposium on the Interface Interface Foundation* (1991), 156–163.

<sup>&</sup>lt;sup>[45]</sup> Rathore, N. et al. The Journal of Chemical Physics **122**, 024111+ (2005).

<sup>&</sup>lt;sup>[46]</sup> Kofke, D. A. The Journal of Chemical Physics 117, 6911–6914 (2002).

że prawdopodobieństwo zamiany replik  $P(T_1 \rightarrow T_2)$  pomiędzy temperaturami T<sub>1</sub> i T<sub>2</sub> zależy od nakrywania się gęstości stanów w tych temperaturach. Prawdopodobieństwo to można obliczyć, o ile znana jest funkcja gęstości stanów  $\Omega(\mathcal{E})$ . Ustalanie temperatur  $T_i$  rozpoczyna się więc od symulacji REMC w której budowane jest pierwsze przybliżenie do  $\Omega(\mathcal{E})$  badanego układu. Na podstawie tej gęstości stanów oblicza się (poprzez numeryczne całkowanie) taki zestaw temperatur, w którym prawdopodobieństwa  $P(T_i \rightarrow T_{i+1})$  były równe dla każdego *i*. Można pokazać, że takie kryterium gwarantuje najszybszą podróż replik pomiędzy temperaturami.

## Inżynieria białek in silico

Najnowszym zastosowaniem pakietu BioShell jest racjonalna inżynieria białek. W pracy H12. opisano udaną modyfikację aldolazy treoninowej z bakterii Thermotoga maritima mającą na celu zwiększenie stabilności tego enzymu. Białko to w żywych organizmach rozkłada treoninę do glicyny i metanalu i w formie aktywnej jest homotetramerem. Dlatego w pracy H12. za cel obrano wzmocnienie oddziaływań pomiędzy łańcuchami. Pierwszym etapem części teoretycznej projektu było zgromadzenie wszystkich znanych sekwencji homologicznych. Wykorzystano do tego program PsiBlastSearch, który uruchamia obliczenia programem PsiBlast z różnymi ustawieniami. Wyniki przeanalizowano narzędziem PsiBlastAnalyse za pomocą którego wybrano 132 reprezentatywne sekwencje. Sekwencje te wykorzystano jako zapytania do kolejnej serii poszukiwań w bazie danych. Ostatecznie znaleziono ponad 100 000 podobnych sekwencji w tym około 45% w drugiej rundzie. 52 z tych sekwencji pochodziło z organizmów termofilnych. Na ich podstawie stworzono uliniowienie wielu sekwencji, pokazujące zmienność aminokwasową na każdej pozycji w białku. Program StrCalc wykorzystano do przeanalizowania struktury krystalograficznej wyjściowego enzymu (depozyt PDB: 1LW5). Na podstawie odległości między atomami i wzajemnej orientacji przestrzennej łańcuchów bocznych reszt aminokwasowych wytypowano potencjalne miejsca do wprowadzenie nowych oddziaływań: mostków jonowych i wiązań disulfidowych.

Struktura czwartorzędowa aldolazy treoninowej jest dość szczególna: pary łańcuchów B i C oraz A i D tetrameru kontaktują się resztami o tych samych numerach. Dla przykładu, prolina 56 z łańcucha A jest oddalona tylko o 4.5Å od P56 w łańcuchu D. Dzięki temu możliwe było wprowadzenie do tetrameru aż czterech reszt cysteiny i jednocześnie *dwóch* wiązań disulfidowych za pomocą jednej tylko mutacji. Po wstępnej analizie do wprowadzania mutacji wybrano 18 pozycji. Modele struktur mutantów obliczono programami Modeller i Rosetta. Ostatecznie 10 najlepszych mutantów przetestowano eksperymentalnie. Dwa z nich (P56C oraz A21C) wykazują znacząco większą stabilność niż białko dzikie.

### 5 Omówienie pozostałych osiągnięć naukowo-badawczych

Habilitant jest również jednym z dwudziestu siedmiu *Principal Investigators* zrzeszonych w *Rosetta Commons*<sup>13</sup> i bierze czynny udział w rozwoju pakietu oprogramowania naukowego Rosetta. Na pakiet ten składa się wiele programów, liczących obecnie w sumie ponad 2,5 miliona linii kodu źródłowego napisanego w języku C<sup>++</sup>. Oprogramowanie to służy do modelowania struktur białek i RNA (zarówno *de novo* jak i w oparciu o szablony), rozwiązywania struktur biomolekuł z fragmentarycznych danych eksperymentalnych (głównie NMR oraz map gęstości elektronowej EM) oraz do projektowania nowych białek. Główny wkład habilitanta to stworzenie algorytmu do budowania biblioteki fragmentów<sup>[40]</sup>, niezbędnej w przypadku modelowania struktur białek tą metodą. Dodatkowo, zaimplementował on umożliwiające wykorzystanie danych SAXS w modelowaniu.

Dominh Con

<sup>&</sup>lt;sup>13</sup> https://www.rosettacommons.org/