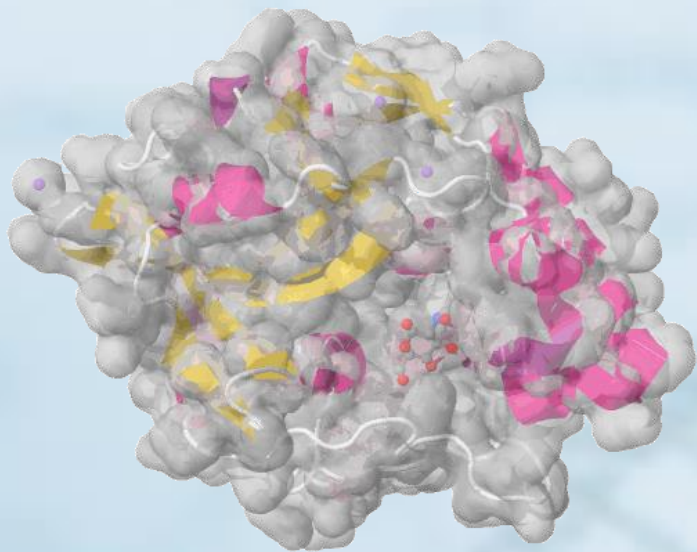


Využití metod strojového učení v bioinformatice

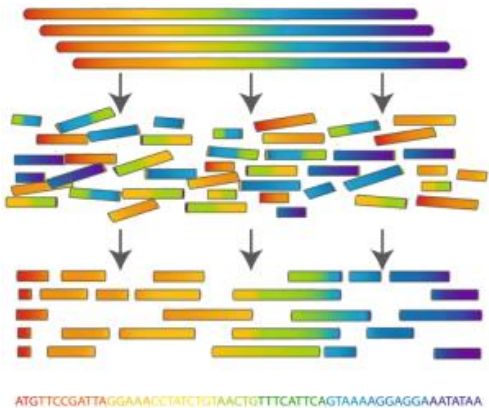


David Hoksza

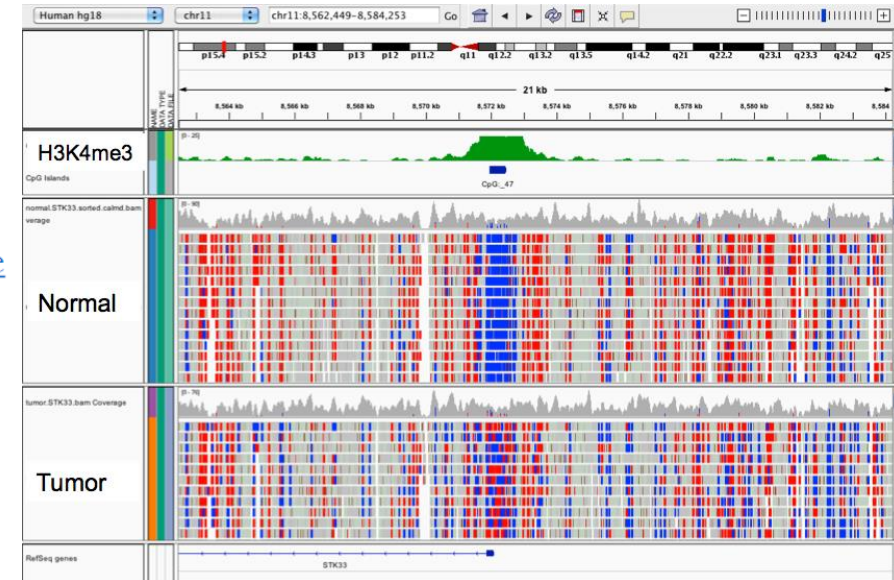
SIRET Research Group
Katedra softwarového inženýrství,
Matematicko-fyzikální fakulta
Karlova Univerzita v Praze

Bioinformatika

- Obor zabývající se vývojem a aplikací počítačových metod pro získávání, ukládání, analýzu a interpretaci biologických dat
 - Zaměření na molekuly DNA, RNA a proteinů



- [cAMP- and cGMP-dependent protein kinase phosphorylation site](#)
 - [RK] (2) -x- [ST]
- [Tyrosine kinase phosphorylation site](#)
 - [RK] -x (2) - [DE] -x (3) -Y or [RK] -x (3) - [DE] -x (2) -Y



Typické úlohy v bioinformatice

- Skládání a mapování genomů
- Podobnostní modelování sekvencí DNA, RNA, proteinů
- Fylogenetická analýza
- Analýza DNA/RNA
 - Hledání sekvenčních motivů
 - Hledání mutací ve vybraných populacích
- ...

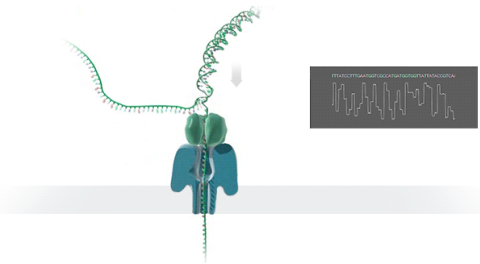
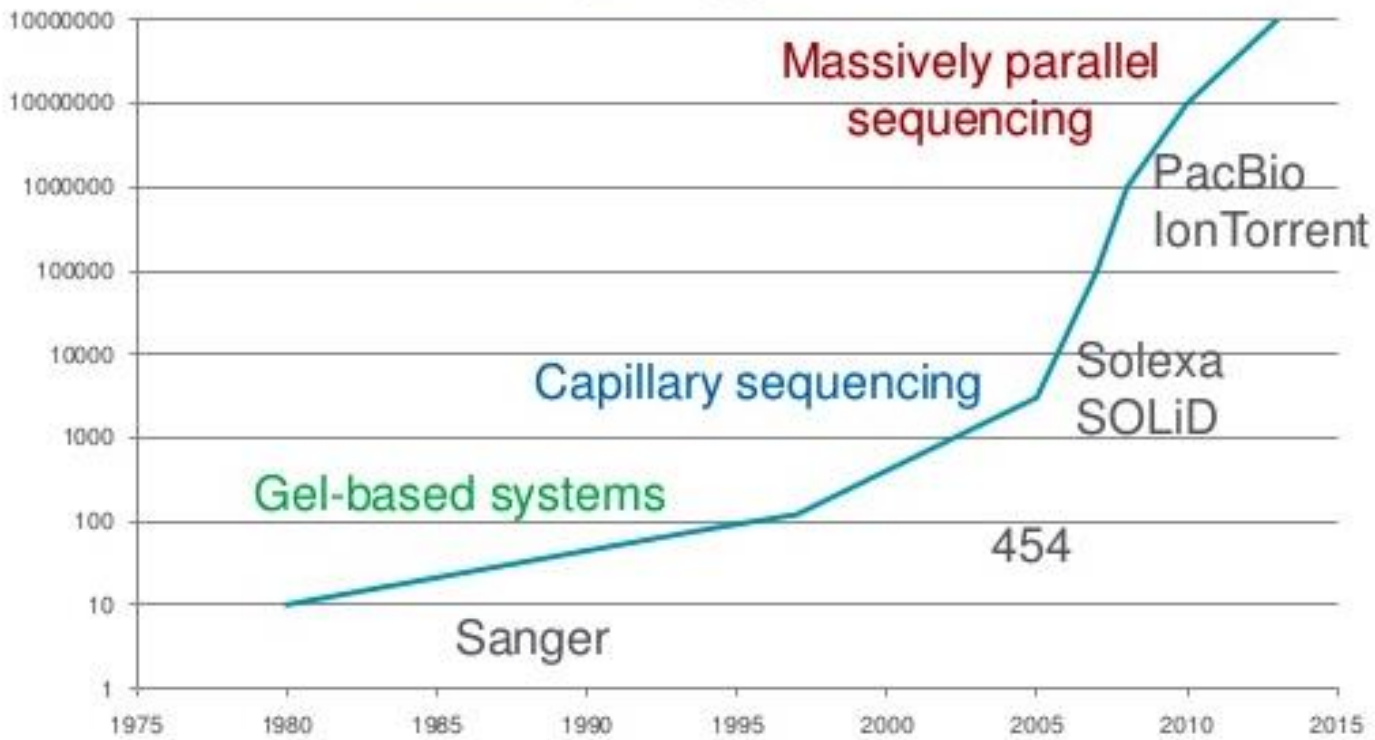
- Predikce struktury proteinů (RNA)
- Identifikace interakcí
- Molekulární modelování (docking)
- ...

- Data genové exprese
- Data hmotnostní spektrometrie
- ...

Aplikační domény bioinformatiky

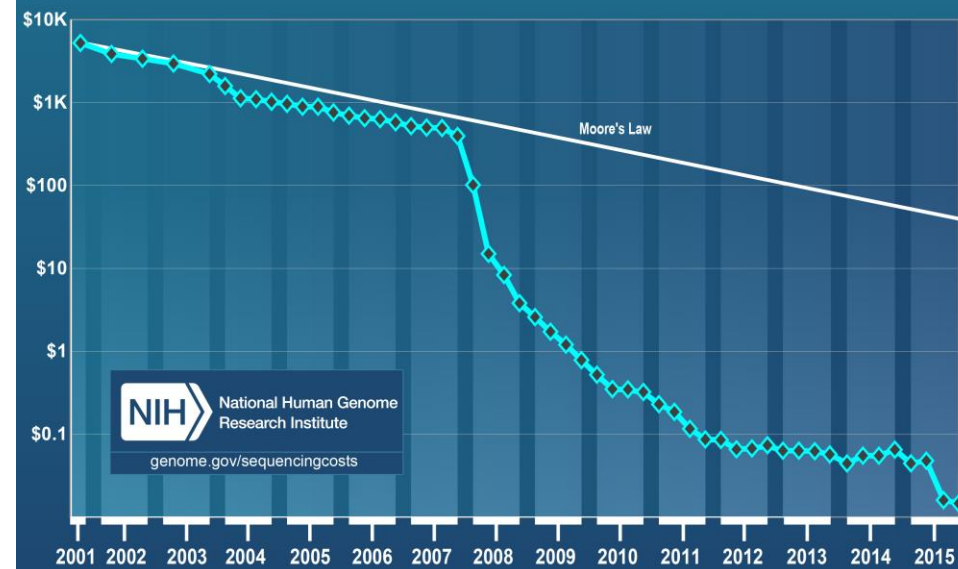
- **Základní výzkum**
 - Funkční genomika, evoluční genomika
- **Medicína**
 - Genetické komponenty nemocí, preventivní medicína, personalizovaná medicína, genová terapie, počítačový vývoj léčiv
- **Mikrobiologie**
 - Biotechnologie, epidemiologie, boj se znečištěním
- **Zemědělství**
 - Odolnější plodiny, resistance proti škůdcům, zvýšení nutriční hodnoty plodin

Kilobases per day per machine

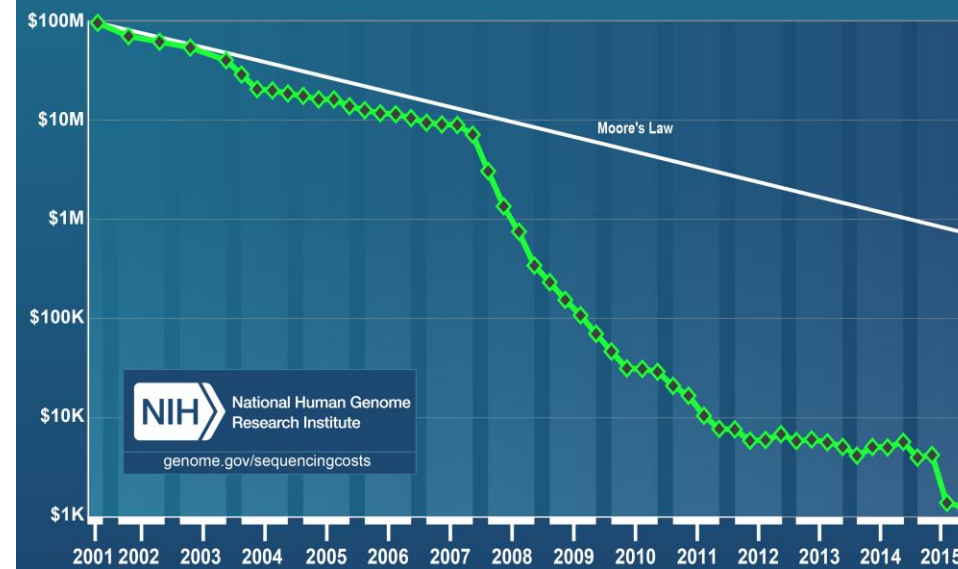


Biologické inspirace informatiky (6. 4. 2017)

Cost per Raw Megabase of DNA Sequence



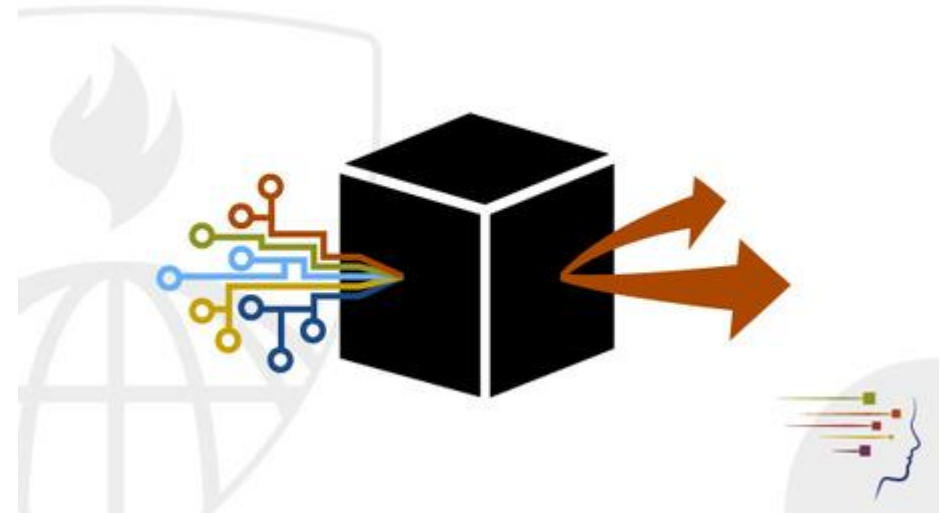
Cost per Genome



Strojové učení

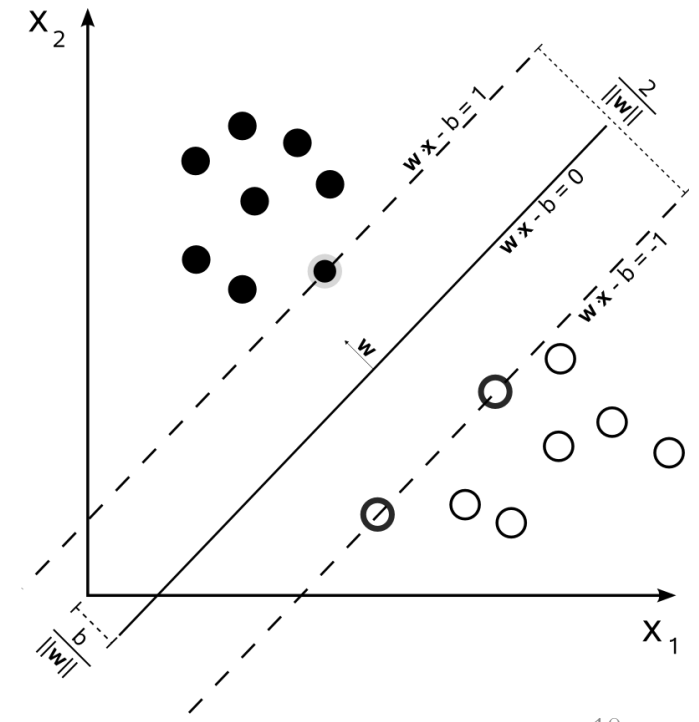
- Množství dat
- Počítačový výkon je dostupný a levný
- Na SU lze nahlížet jako na jiný způsob programování

} Rozhodovací procesy,
predikce



Strojové učení

- Skupina algoritmů schopných **identifikovat (naučit se) vzory v datech** a tuto znalost **aplikovat na dosud neviděných příkladech**
 - Regrese
 - **Klasifikace**
 - Shlukování
- Typy
 - **Strojové učení s učitelem (supervised learning)**
 - Strojové učení bez učitele (unsupervised learning)
- Typické přístupy
 - Rozhodovací stromy
 - Support vector machines
 - Neuronové sítě
 - Markovovské modely
 - Pravděpodobnostní grafické modely



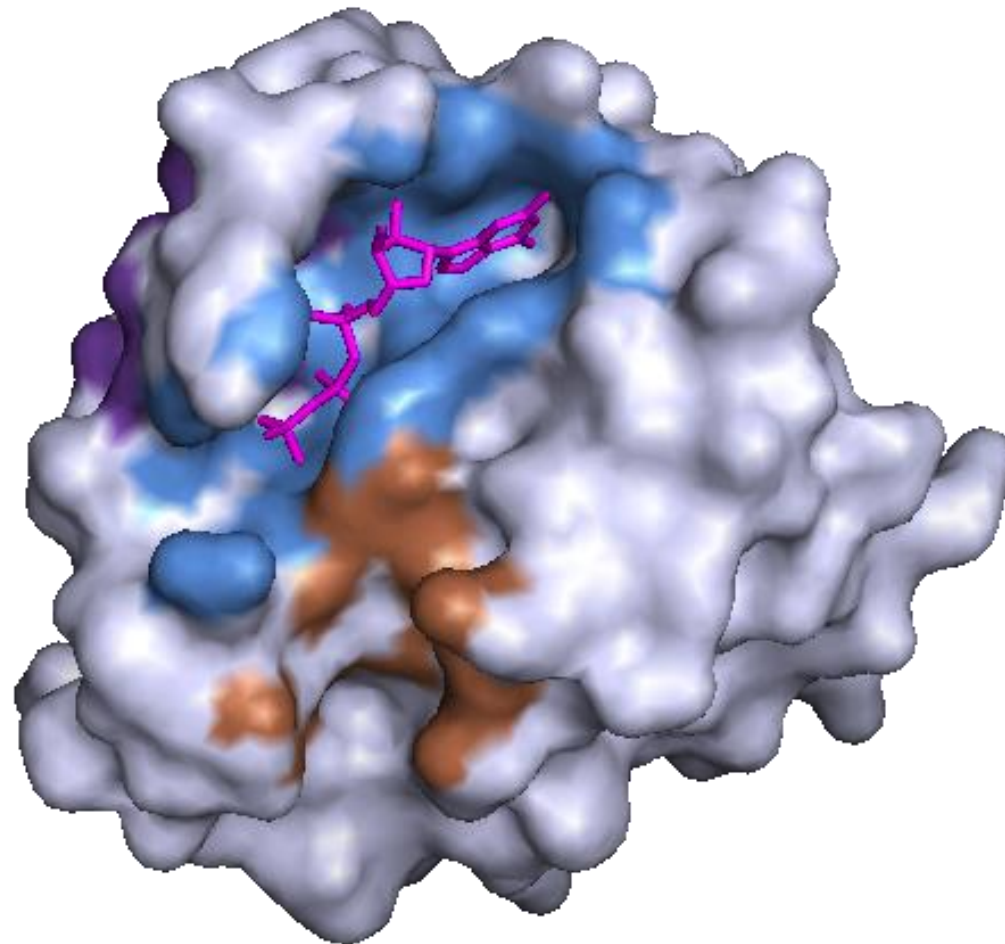
Bioinformatické úlohy řešené strojovým učením

- Identifikace sekvenčních motivů
 - Rozpoznávání a funkční anotace genů
 - Hledání protein-vázajících míst v DNA
 - Rozpoznání sekundární struktury
 - Opravy sekvenačních chyb
- Identifikace signálů
- Predikce vazebných míst na úrovni struktury
- Predikce bioaktivních malých molekul (léčiv)
- ...

Identifikace (protein-ligand) aktivních míst na povrchu proteinu

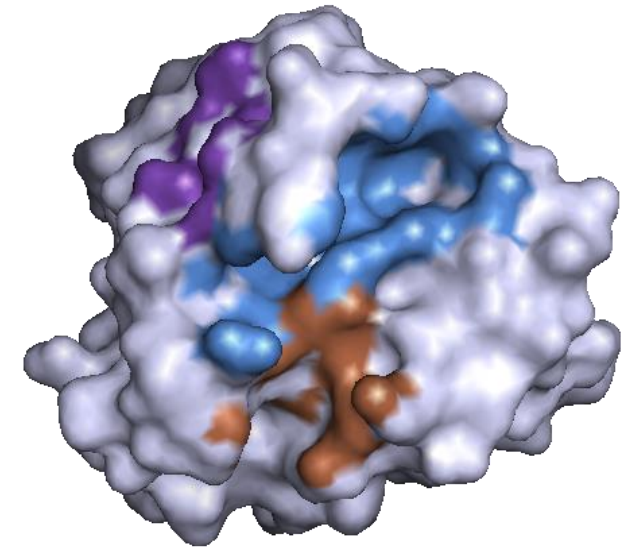
Protein-ligand interakce

- **Aktivní místa (kapsy)**
- **Motivace**
 - Identifikace potenciálních **cílů pro léčiva**
 - Predikce **vedlejších účinků** léčiv
 - **Predikce funkce** neznámého proteinu



P2RANK

- Počítačová metoda (algoritmus) schopný **identifikovat místa na povrchu proteinu**, na které se může s vysokou pravděpodobností **vázat nspecifikovaný ligand**
- **Vstup:** počítačová reprezentace proteinové struktury
- **Výstup:** seznam míst na povrchu proteinu pravděpodobně schopných vázat ligand



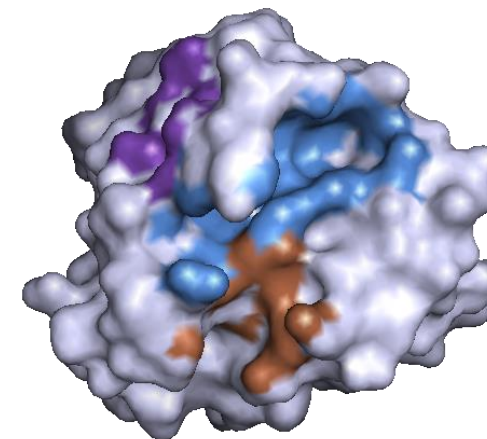
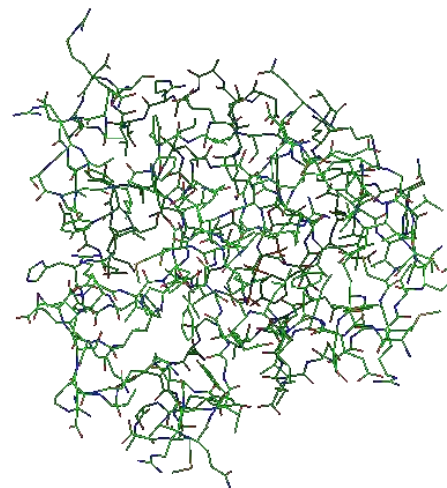
Informatický pohled na protein

Sekvence

- **Řetěz** aminokyselin → **lineární** **sekvence** písmen (slovo)
- Písmena reprezentují aminokyseliny (ARNDCCEQGHILKMFPSTWYV)

Struktura

- **Pozice** jednotlivých atomů v 3D prostoru

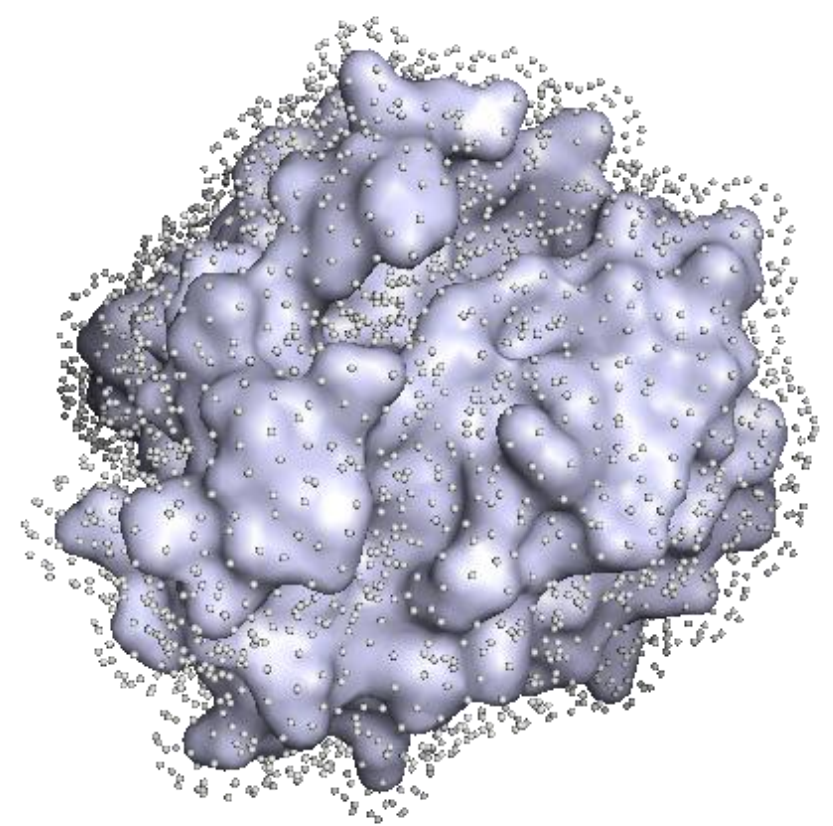


P2RANK princip

- **Využití informací o existujících** aktivních místech a jejich rysech pro **rozpoznání typově podobných** míst na neznámém proteinu
 - **Jak určit kapsu na dosud neviděném proteinu?**
 - Strojové učení (s učitelem)
 - Fáze učení: naučení modelu pro rozpoznání rysů bodů aktivních míst
 - Fáze rozpoznávání: aplikace modelu na povrch neznámého proteinu
 - **Jak popsat rysy povrchu proteinu?**
 - Projekce fyzikálně chemických vlastností aminokyselin na povrch proteinu

Algoritmus – učící fáze

1. Získání známých protein-ligand komplexů
2. Potažení povrchů proteinů **sítí bodů**
3. **Extrakce vektoru fyzikálně-chemických vlastností** pro každý z bodů každého proteinu
4. **Vybudování modelu**, který pro daný bod (vektor) bude schopný určit, s jakou pravděpodobností je tento součástí kapsy = **strojové učení**

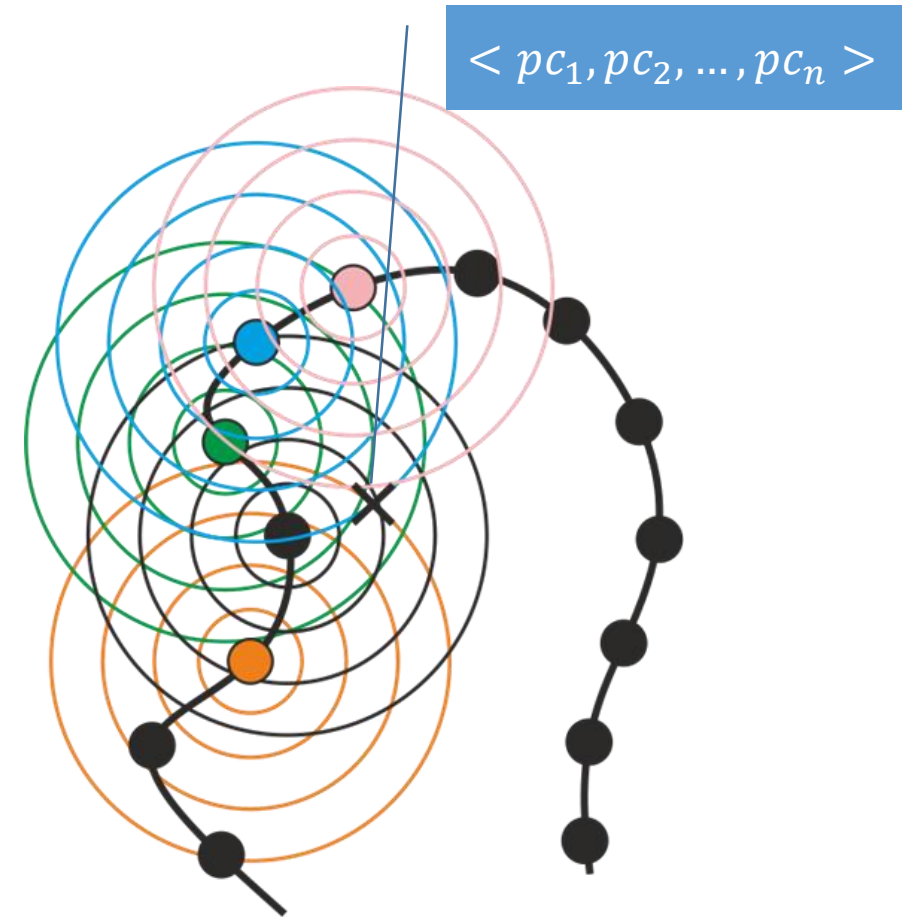


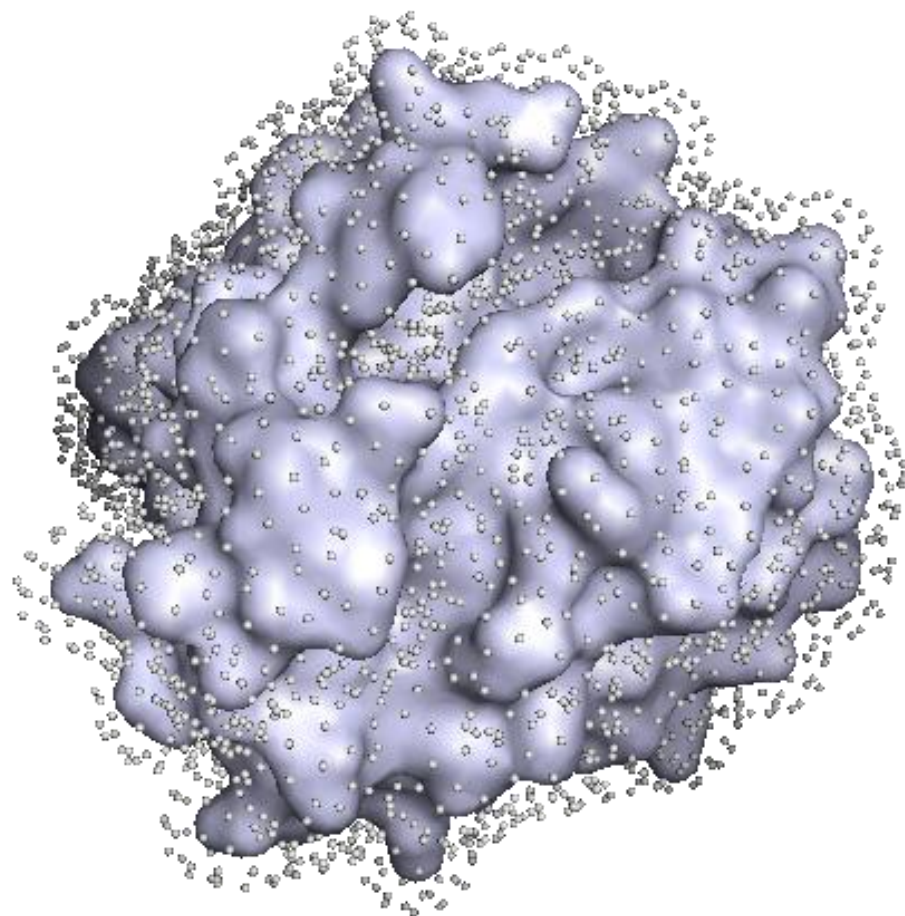
P2RANK - extrakce vlastností

- Okolo **30 atributů** popisující fyzikálně-chemické vlastnosti aminokyselin a lokálního okolí daného bodu

$$PC(V) = \frac{1}{m} \sum_{A_i \in A(V)}^m PC(A_i) \cdot w(dist(V, A_i))$$

$$w(d) = \begin{cases} 1, & d \leq 4 \text{ \AA} \\ (4/d)^2 & d > 4 \text{ \AA} \end{cases}$$

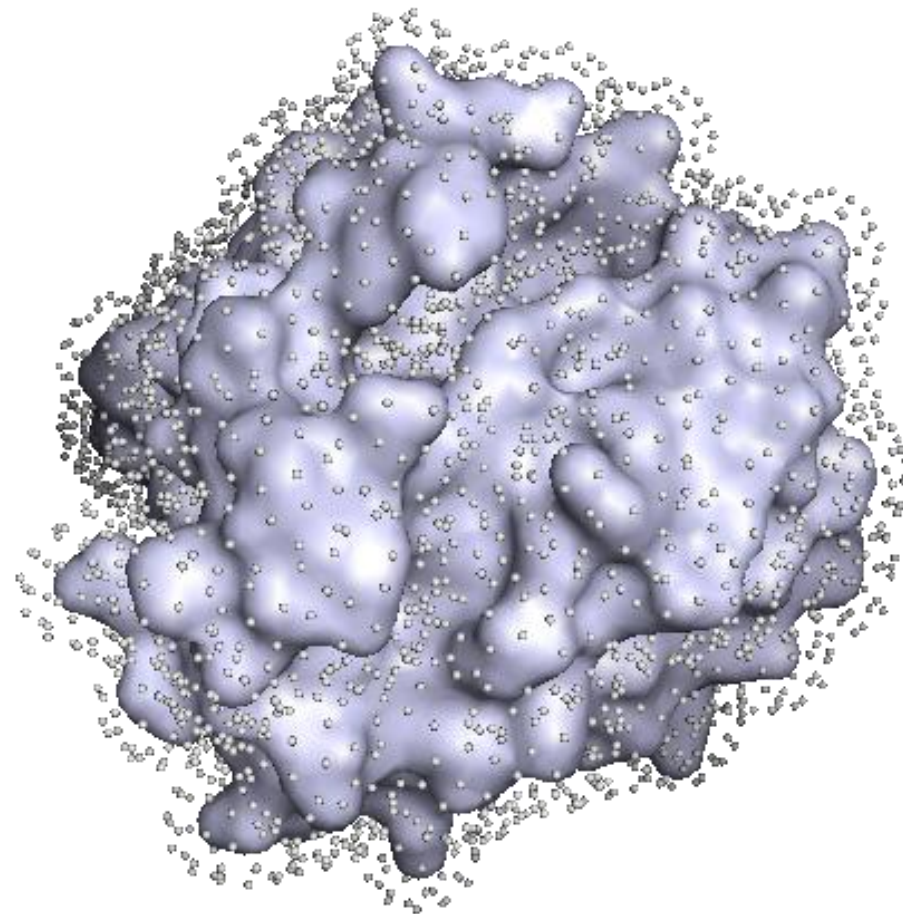
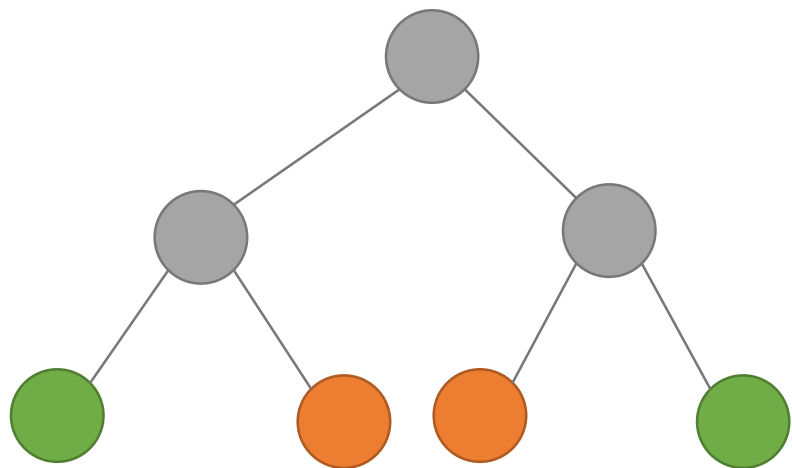




Rozhodovací stromy

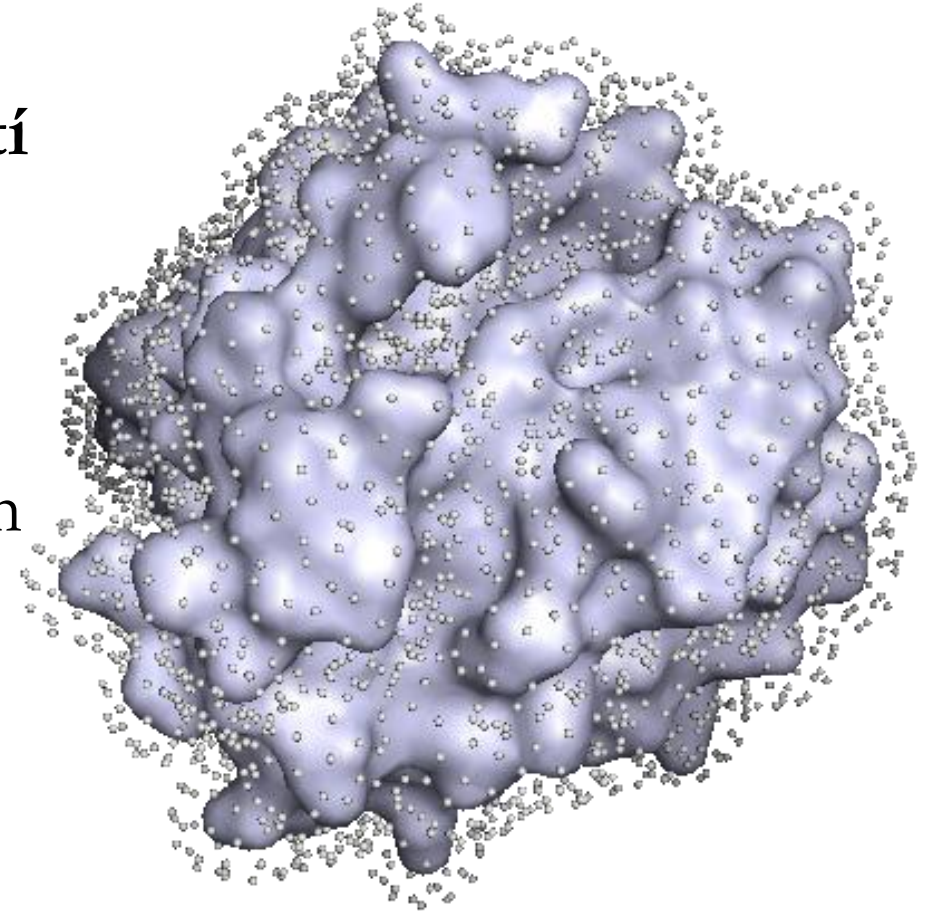
- Model reprezentován **stromem**

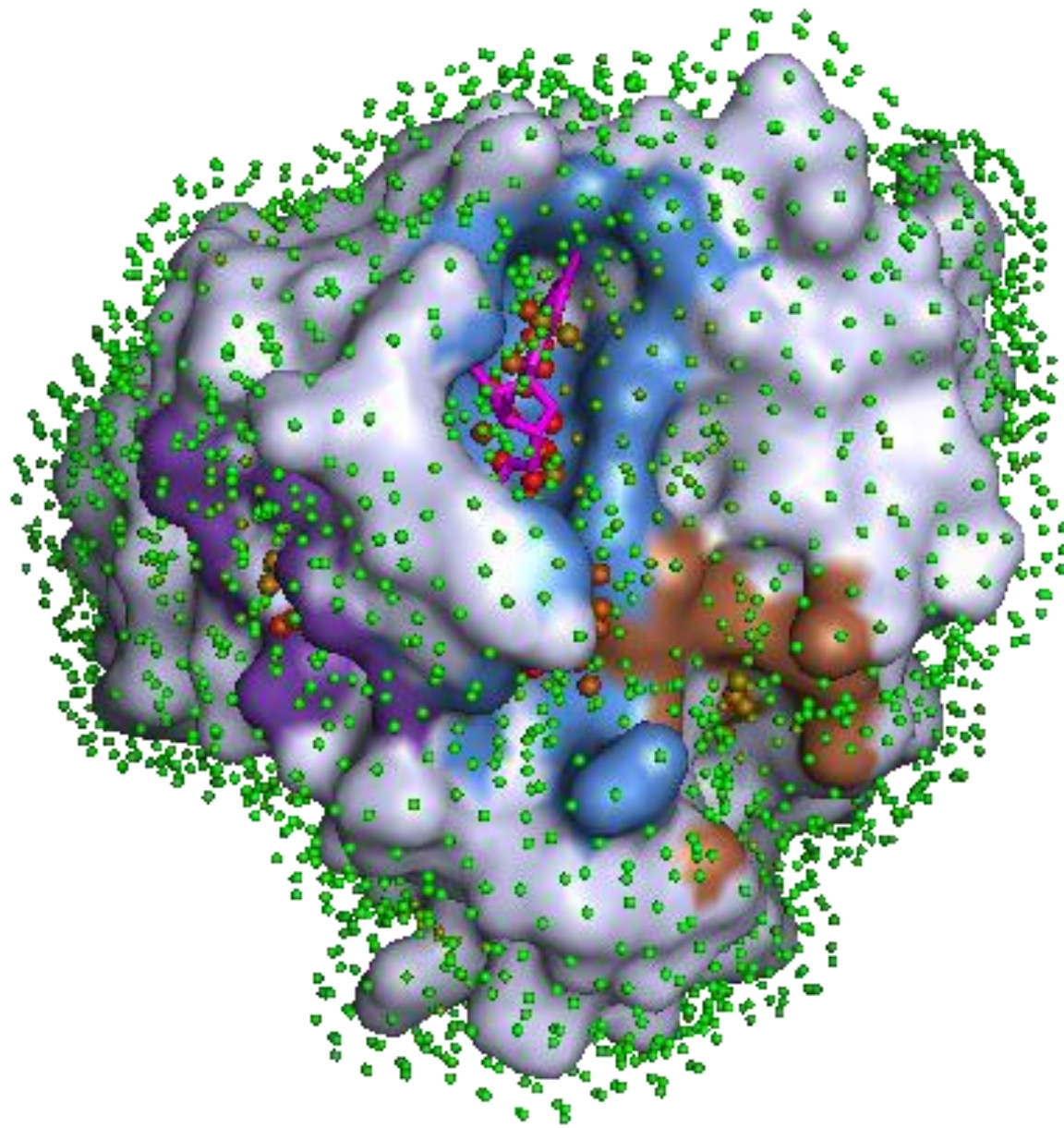
- Vnitřní uzly pravidla/otázky
- Hrany odpovědi
- Listy rozhodnutí



Algoritmus – fáze rozpoznávání

1. **Potažení povrchu** neznámého proteinu **sítí bodů**
2. **Aplikace modelu** pro každý bod sítě → vazebné skóre bodu
3. Vypuštění bodů s nízkým vazebným skórem
4. **Identifikace shluků** vysoce skórujících bodů → **kapsa**
5. **Ohodnocení** kapes součtem skór jejich bodů





Prank Download Authors Help

Root Entity

- ▼ **M** PDB
- ▼ **M_M** Model 1 1225 atoms
 - ▼ **S_M** Protein complement
 - ⊕ **V_M** Surface, 0.5 Å probe
 - ▼ **S_M** pocket1 23 atoms
 - ⊕ **V_M** Balls and Sticks
 - ⊕ **V_M** Surface, 0.5 Å probe
 - ▼ **S_M** pocket2 24 atoms
 - ⊕ **V_M** Balls and Sticks
 - ⊕ **V_M** Surface, 0.5 Å probe
 - ▼ **S_M** pocket3 21 atoms
 - ⊕ **V_M** Balls and Sticks
 - ⊕ **V_M** Surface, 0.5 Å probe
 - ⊕ **V_M** Surface, 0.7 Å probe
 - ▼ **S_M** pocket4 17 atoms
 - ⊕ **V_M** Balls and Sticks
 - ⊕ **V_M** Surface, 0.5 Å probe
 - ▼ **S_M** pocket5 14 atoms
 - ⊕ **V_M** Balls and Sticks
 - ⊕ **V_M** Surface, 0.5 Å probe
- ▼ **G** Group Macromolecule
 - ▼ **S_M** Polymer 1060 atoms
 - ⊕ **V_M** Cartoon
 - ▼ **S_M** HET 26 atoms
 - ⊕ **V_M** Balls and Sticks

Pockets

50 100

Model 1 1225 atoms

- S_M** Selection
- M_M** Crystal Symmetry
- A** Macromolecule Visual
- ✓ Add
- V_M** Visual
 - + Type Cartoon
 - + Coloring Chain ID
 - ✓ Add
- M_M** Update Model

Pockets:

⊕	⊖
Pocket rank: 1	Pocket score: 8.0031
AA count: 10	Conservation: N/A
⊕	⊖
Pocket rank: 2	Pocket score: 2.4908
AA count: 13	Conservation: N/A
⊕	⊖
Pocket rank: 3	Pocket score: 2.0346
AA count: 9	Conservation: N/A
⊕	⊖
Pocket rank: 4	Pocket score: 1.1163
AA count: 7	Conservation: N/A
⊕	⊖
Pocket rank: 5	Pocket score: 0.8559
AA count: 9	Conservation: N/A

Jaké znalosti jsou třeba pro vývoj P2RANKu?

- Základy biologie, proteomiky
- Znalost zdrojů biologických dat (PDB)
- Programování
- Pokročilá algoritmizace
- Pravděpodobnost a statistika


Bioinformatika na Univerzitě Karlově v Praze

Zdroj informací o bioinformatice nejen na UK

<http://bioinformatika.mff.cuni.cz/>

 Studijní program

 Svět bioinformatiky

 Bioinformatický seminář

Bioinformatika je obor, ve kterém se setkává biologie a informatika při studiu biologických dat jako jsou DNA, RNA a proteiny.

Biologie těží ve 21. století z neobyčejného rozvoje metodických přístupů (sekvenování, proteomika), který vedlo a dále vede k bezprecedentnímu nárůstu cenných biologických dat. Je velmi pravděpodobné, že se v dohledné době stane sekvenování kompletního genomu standardním diagnostickým vyšetřením a Evropan absolvuje sekvenování svého genomu několikrát za život. I

Dotazy

