

# *Bioinformatická analýza mutací*

Kamila Réblová



## Rozvoj nástrojů pro data mining:

- rozvoj bioinformatických nástrojů (např. strukturní a predikční metody)
- vznik odvozených databází



## Archivace sekvenčních dat (databáze), snadná manipulace, dostupnost

Primární (sekvenční) databáze nukleových kyselin:

NCBI (National Center for Biotechnology Information)

-> GenBank <http://www.ncbi.nlm.nih.gov/Genbank/>

EBI (European Bioinformatics Institute) -> European Molecular Biology Laboratory (EMBL)-Bank <http://www.ebi.ac.uk/embl/>

DNA Data Bank of Japan (DDBJ)

<https://www.ddbj.nig.ac.jp/index-e.html>

# Formáty sekvencí:

## Fasta formát

```
>sp|P00439|PH4H_HUMAN Phenylalanine-4-hydroxylase OS=Homo sapiens GN=PAH PE=1 SV=1
MSTAVLENPGLGRKLSDFGQETSYIEDNCNQNNGAISLIFSLKEEVGALAKVLRLFEENDV NLTHIESRPSRLKKD
EYEFFTHLDKRSPLALTNIILRHDIGATVHELSDKDKKDTVPW FPRTIQELDRFANQILSYGAELDADHPGFKD
PVYRARRKQFADIAYNYRHGQPIPRVEYM EEEKKTWGTVFKTLKSLYKTHACYEYNHIFPLLEKYCGFHEDNIP
QLEDVSQFLQTCTGF RLRPVAGLLSSRDFLGGLAFRVFHCTQYIRHGSKPMYTPEDICHELLGHVPLFSDRSF
A QFSQEIGLASLGAPDEYIEKLATIWFTVEFGLCKQGDSIKAYGAGLLSSFGEQYCLSE KPKLLPLELEKTAIQ
NYTVTEFQPLYVAESFNDAKEKVRNFAATIPRPFSVRYDPYTQR IEVLDNTQQLKILADSINSEIGILCSALQKIK
```

## GenBank format

A sequence file in GenBank format can contain several sequences.  
One sequence in GenBank format starts with a line containing the word LOCUS and a number of annotation lines.

An example sequence in GenBank format is:

```
LOCUS      AB000263                368 bp    mRNA    linear    PRI 05-FEB-1999
DEFINITION Homo sapiens mRNA for prepro cortistatin like peptide, complete
            cds.
ACCESSION  AB000263
ORIGIN
    1  acaagatgcc attgtccccc ggcctcctgc tgetgetget ctccggggcc acggccaccg
   61  ctgccctgcc cctggagggg ggcctccacc gccgagacag cgagcatatg caggaagcgg
  121  caggaataag gaaaagcagc ctctgactt tctctgcttg gtggtttgag tggacctccc
  181  aggccagtgc cgggcccctc ataggagagg aagctcggga ggtggccagg cggcaggaag
  241  gcgcaccccc ccagcaatcc gcgcgcgggg acagaatgcc ctgcaggaac ttctctgga
  301  agaccttctc ctctgcaaa taaaacctca cccatgaatg ctacgcgaag ttttaattaca
  361  gacctgaa
```

//

---

# IUPAC nucleotide code

## Nucleotide ambiguity code

(as defined in DNA Baser Sequence Assembler)

Code	Represents	Complement
A	Adenine	T
G	Guanine	C
C	Cytosine	G
T	Thymine	A
Y	Pyrimidine (C or T)	R
R	Purine (A or G)	Y
W	weak (A or T)	W
S	strong (G or C)	S
K	keto (T or G)	M
M	amino (C or A)	K
D	A, G, T (not C)	H
V	A, C, G (not T)	B
H	A, C, T (not G)	D
B	C, G, T (not A)	V
X/N	any base	X/N
-	Gap	-

## Large-scale DNA sequencing projects:

1953 - Odhalena struktura DNA... ???

1965 – přečtena sekvence tRNA kvasinky

1974 – přečtena sekvence genomu bakteriofága  $\Phi$ X174 (5,375 bp)

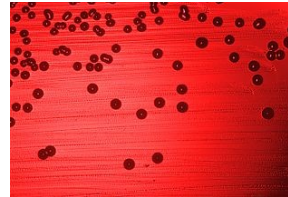
1988-1990 – založen Human Genome Project (HGP)

(probíhá 1990-2000/2003)

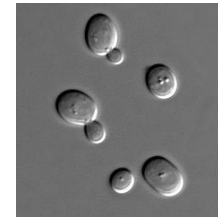
### První organismy:

1995 – *Haemophilus influenzae* (1,830,140 bp)

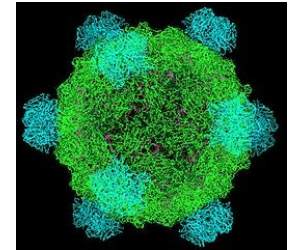
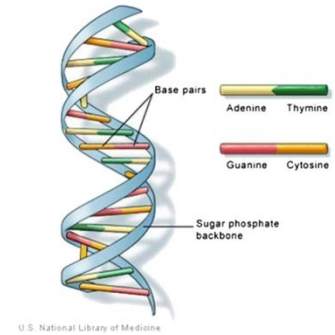
(*epiglottitida, meningitida*)



1996 – *Saccharomyces cerevisiae* (12,068,000 bp)



1998 – *Caenorhabditis elegans* (100,000,000 bp)



## Large-scale DNA sequencing projects:

**1953 - Odhalena struktura DNA, J. Watson, F. Crick, R. Franklin**

**1965 – přečtena sekvence tRNA kvasinky**

**1974 – přečtena sekvence genomu bakteriofága  $\Phi$ X174 (5,375 bp)**

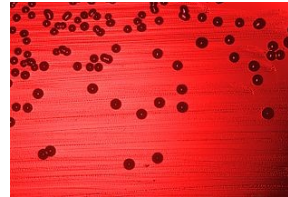
**1988-1990 – založen Human Genome Project (HGP)**

(probíhá 1990-2000/2003)

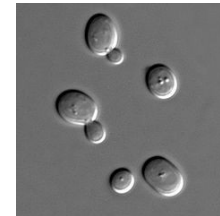
### První organismy:

**1995 – Haemophilus influenzae (1,830,140 bp)**

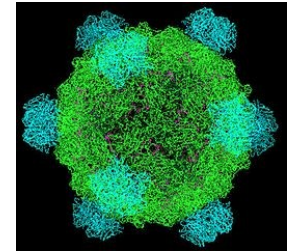
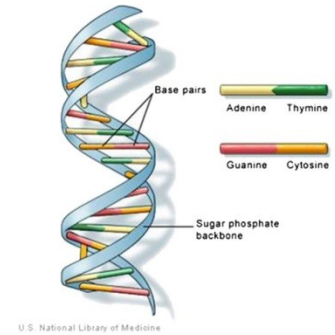
(*epiglottitida, meningitida*)



**1996 – Saccharomyces cerevisiae (12,068,000 bp)**



**1998 – Caenorhabditis elegans (100,000,000 bp)**

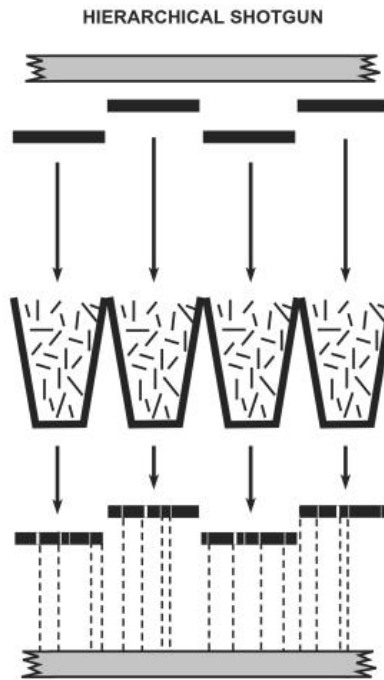
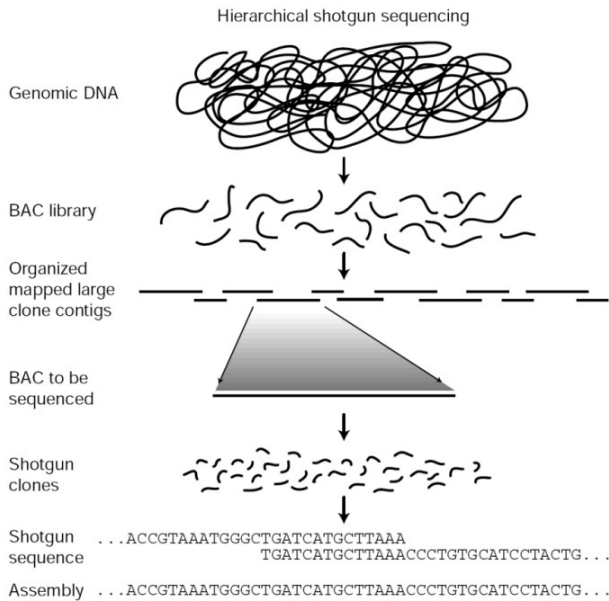


2001 – Lidský genom (3,200,000,000 bp)

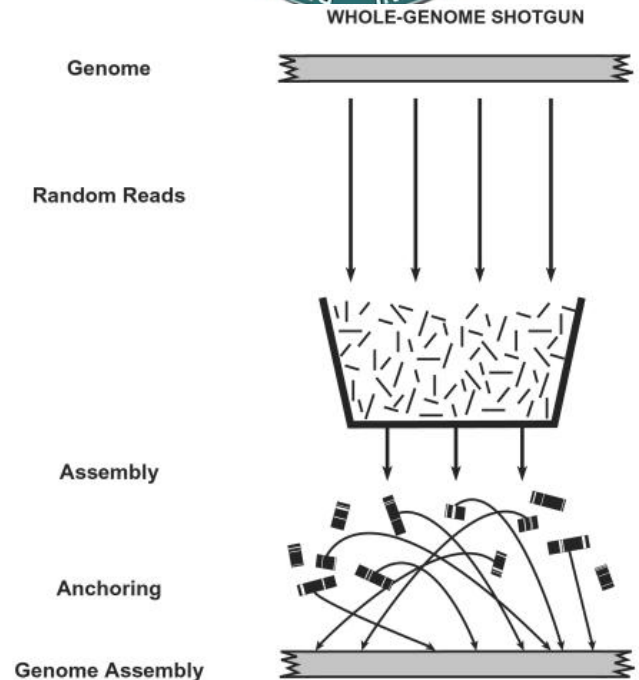
Human Genome Project Consortium (2001) "A physical map of the human genome" Nature **409**:934–41

Dideoxy neboli 'Sangerova metoda'

1977 – zavedena Sangerova metoda  
(Nobelova cena 1980)



HGP



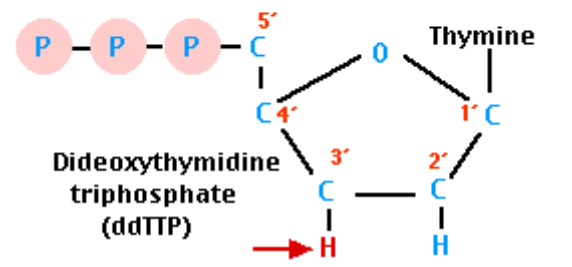
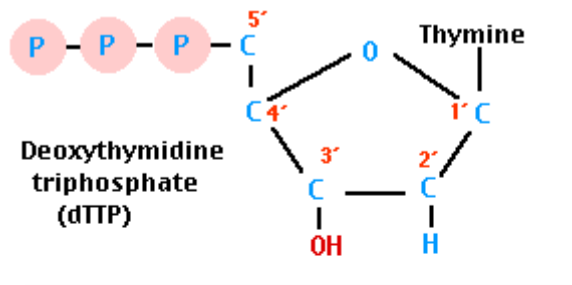
Celera Company



# Sangerova metoda

Původně gelová elektroforéza a radioaktivní značení  
=> fluorescenčních značení

**dideoxy**nukleotidtrifosfátů v kombinaci s kapilární elektroforézou.  
Takto značené oligonukleotidy jsou detekovány laserem

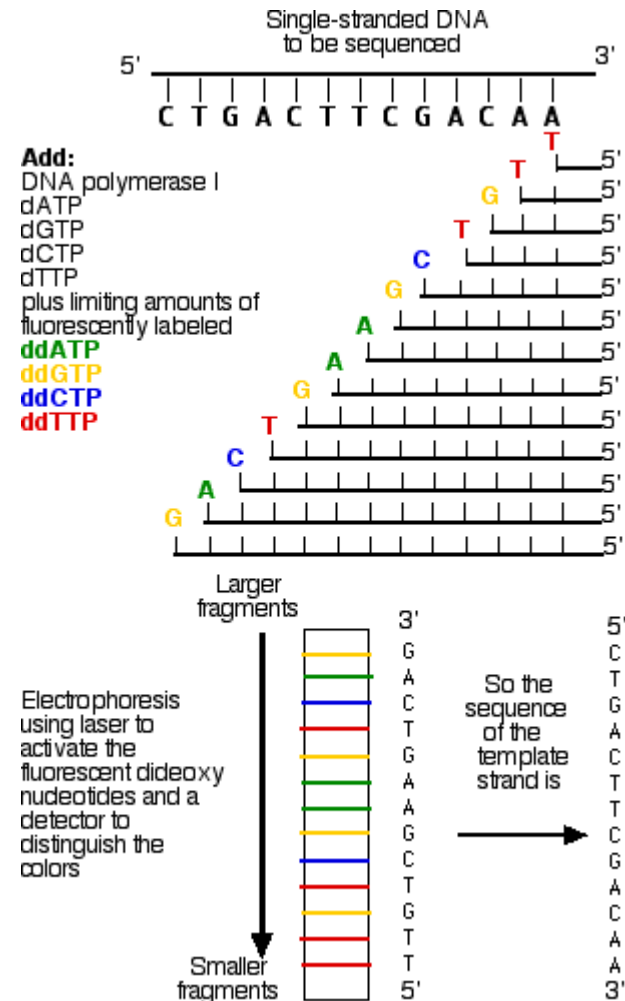


dATP  
dGTP  
dCTP  
dTTP

DNA polymerase I

ddATP  
ddGTP  
ddCTP  
ddTTP

e.g. The concentration of ddATP should be 1% of the concentration of dATP.



3,2 mld párů bází



Genom

2% genomu → proteiny  
zbytek – “junk” DNA -  
Non-coding = tRNA, rRNA,  
Transpozony

Co ovlivňují? - 3D strukturu DNA -  
Chromosomální uspořádání,  
regulace procesů v buňky (dělení)

Proteom

Proteom je soubor proteinů, které jsou produkované z genomu daného organismu  
Je dynamický, závisí na věku, ale také typu tkáně a stavu organismu.

Mikrobiom

1:10? 10:1 ( bakterie ca. 1 kg)

3 200 000 000 párů bází  
23 párů chromozómů  
Délka DNA= 1 metr

známé protein-kódující  
geňy

člověk

?

octomilka

?

3 200 000 000 párů bází  
23 párů chromozómů  
Délka DNA= 1 metr

**známé protein-kódující  
geny**

**člověk**

**20-25 tis**

**octomilka**

**13 tis**

3 200 000 000 párů bází  
23 párů chromozómů  
Délka DNA= 1 metr

	<b>člověk</b>	<b>octomilka</b>
<b>známé protein-kódující geny</b>	<b>20-25 tis</b>	<b>13 tis</b>

<b>transkripty</b>	<b>178 191</b>	<b>23 017</b>
--------------------	----------------	---------------

**95-100% lidských genů podléhá sestřihu, min. 2 transkripty.**

# Next (second) generation sequencing – rozvoj od 2005

**'Next' = not Sanger Method**

**454 sequencing** – pyrosequencing

**Illumina** – sequencing by synthesis/bridge amplification

**SOLiD** - Sequencing by Oligonucleotide Ligation and Detection

**Third generation - Oxford nanopore technology (ONT)**

**Aviti**

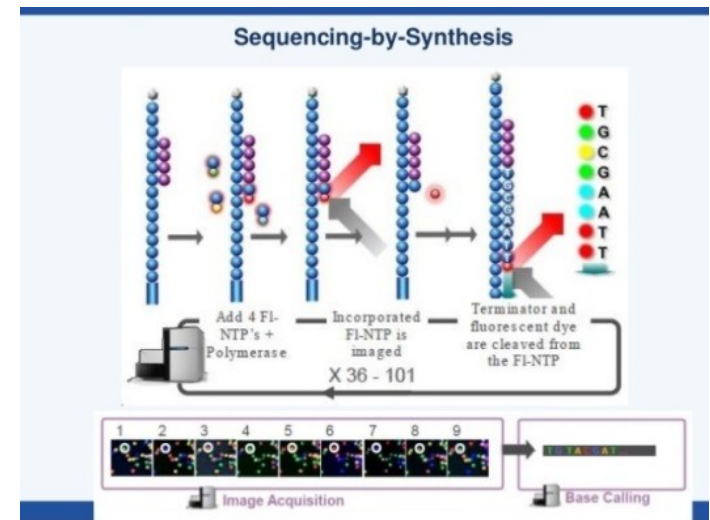
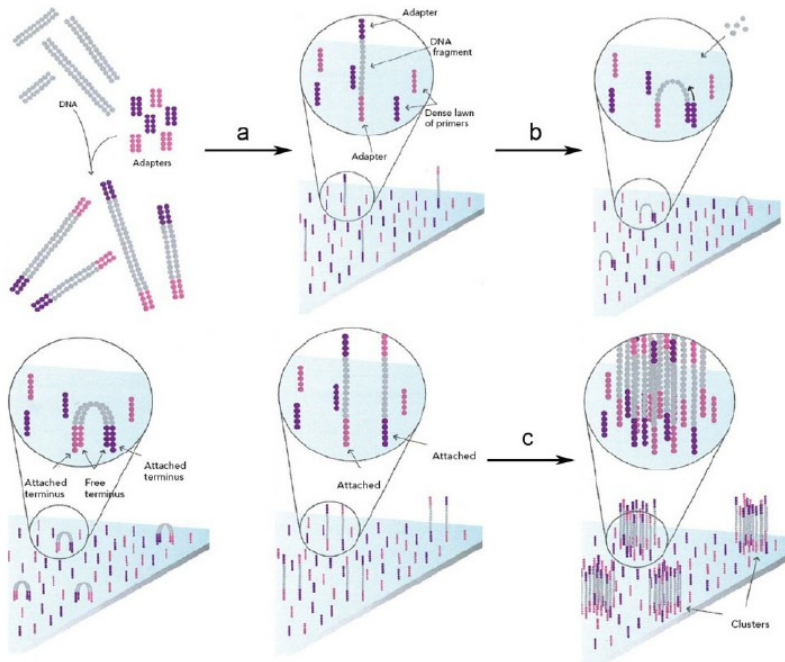
Nové platformy

**MGI**

## Capture seq.

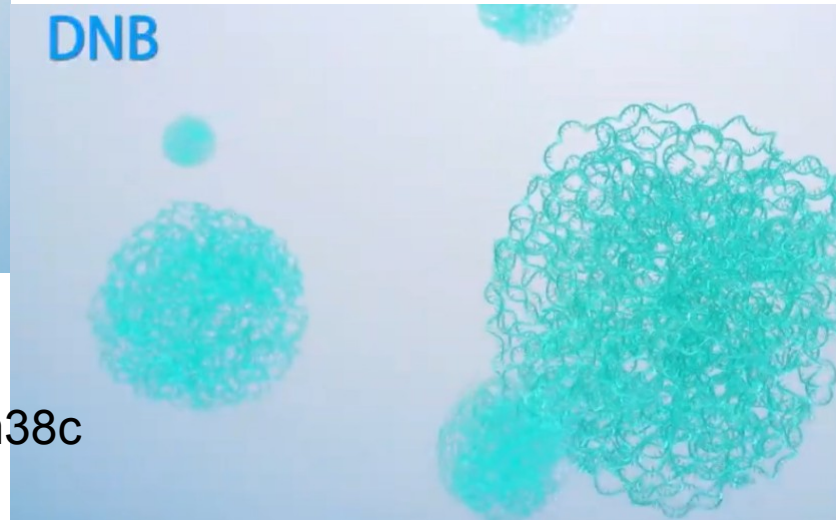
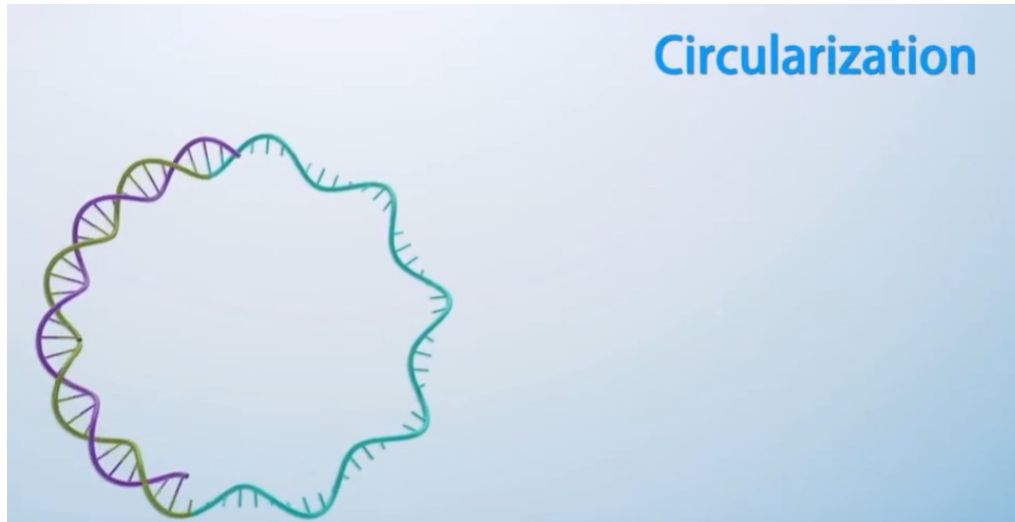
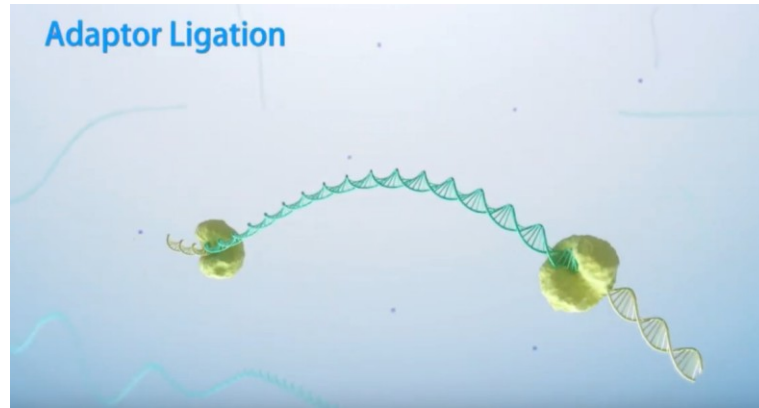
(Použití sond, které vychytají cílené geny - sequence capture)

- Fragmentace (fyzikálně např. sonikace, nebo enzymaticky (DNAázy, Transpozáz) nebo chemicky (tepelné štěpení dvojmocnými ionty))
- Úprava konců – end repair (zarovnání), fosforylace, dATP
- Ligace adaptoru, které představují „primery“
- Vychytání úseků pomocí hybridizačních sond
- Amplifikace



# AVITI a MGI

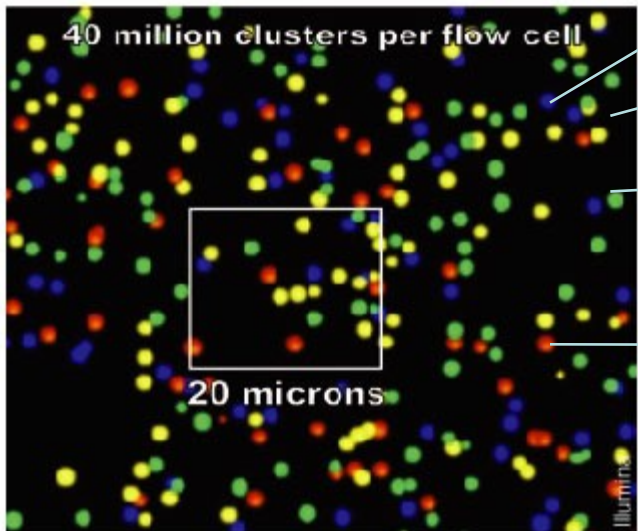
Cirkularizace DNA



<https://www.youtube.com/watch?v=xUVdJN0m38c>



Tif fily -> bcl file (intensity)-> Fastq fily



A

G

C

T

Primární data

Fastq file

```
Mock_A_TATAGCGA-GACACCGT_L001_R1_001.fastq
@HWI-M01141:63:A4NDL:1:1101:16668:1377 1:N:0:TATAGCGAGACACCGT
NACAGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCGCTAGGTGGTTTGT/
+
#>>AA>CAABBBGGGGGGFFGHFEGGGFHHHGHFEGCEHHFEGGGG@EEHHGGGHHC
@HWI-M01141:63:A4NDL:1:1101:14849:1418 1:N:0:TATAGCGAGACACCGT
NACGAGGGTGCAAGCGTTACTCGGAATTACTGGGCGTAAAGCGTGCCTAGGTGGTGGTT/
+
#>>>A??AFAA1BGEGGAAFGGCA0BFF1D2BCF/EEG/DBEE/E?GAEEFGEFAEFG1
@HWI-M01141:63:A4NDL:1:1101:13802:1421 1:N:0:TATAGCGAGACACCGT
NACGAGGGTGCAAGCGTTAATCGGAATTACTGGGCGTAAAGCGCACGCAGGCGTTTGT/
+
#>>AAABBBABBBGGGGGG?FGHGGGGHHHHHHHGGGGHHGGGGGGGGGGGGGG?F1
@HWI-M01141:63:A4NDL:1:1101:15928:1426 1:N:0:TATAGCGAGACACCGT
NACGTAGGTGCGAGCGTTAATCGGAATTACTGGGCGTAAAGCGTGCAGGCGTTTGTG/
+
#>>AABFB@FBBGGGGGGGGGGHGGGGFHHHHHHHGGGGHHGGGGGGGGGGGGEE?G4
@HWI-M01141:63:A4NDL:1:1101:14861:1431 1:N:0:TATAGCGAGACACCGT
NACGAGGGTGCAAGCGTTACTCGGAATTACTGGGCGTAAAGCGTGCCTAGGTGGTGGTT/
+
#>>AAAABFBABGGGGGCEGHGGEFFHHHHHHHGGGGHHGGGGGEGFHHGGGGEGHE
@HWI-M01141:63:A4NDL:1:1101:15264:1465 1:N:0:TATAGCGAGACACCGT
NACGTAGGTGCGAGCGTTGTCCGGAATTACTGGGCGTAAAGAGCTCGTAGGTGGTTTGTCC
+
```

Info o běhu NGS

Sekvence readu

Kvalita readu

Primární sekvence





## 2. Úprava bam souboru

Indexování, realignment indelů, odstranění duplikátů

## 3. Variant calling – statistické testy, porovnává nalezené rozdíly vůči referenci

výsledek: vcf file – soubor nalezených genetických variant

Nastavení pro bioinformatickou proceduru např:

- Threshold for allele frequency – 0,2 = 20% (for germline)
- Mapping quality – 30
- Base quality – phred score
- Počet mismatchů v readu
- aj.

4. Anotace ... u kodujících úseků přeloží do proteinové sekvence.  
 Zjistí četnost dané varianty proti databázím, např: 1000genomů, Clinvar, Gnomad...)  
 Přidá info z Omimu

### Výsledný soubor -> xls file

The screenshot shows a Microsoft Excel spreadsheet titled "JL2589\_AAACAT\_bcf2.hg19\_multianno\_final - Microsoft Excel". The spreadsheet contains a large table of genomic data. The columns are labeled A through BL. The data includes chromosome (CHROM), position (POS), reference (REF), alternative (ALT), quality (QUAL), filter (FILTER), function (Func.refGene), gene (Gene.refGene), and various variant frequencies and annotations from databases like 1000G, Clinvar, and GnomAD. The 'P' column (Clinvar) is highlighted with a yellow background. The 'O' column (hgmd) contains the text 'clinvar'. The 'Q' column (esp6500si) contains 'esp6500si'. The 'R' column (esp6500sl) contains 'esp6500sl'. The 'S' column (JL2589\_AA) contains 'JL2589\_AA'. The 'T' column (JL2589\_AA) contains 'JL2589\_AA'. The 'U' column (JL2589\_AA) contains 'JL2589\_AA'. The 'V' column (JL2589\_AA) contains 'JL2589\_AA'. The 'W' column (JL2589\_AA) contains 'JL2589\_AA'. The 'X' column (ADJAF) contains 'ADJAF'. The 'Y' column (HIAF) contains 'HIAF'. The 'Z' column (QUAL) contains 'QUAL'. The 'AA' column (SBF) contains 'SBF'. The 'AB' column (MQ) contains 'MQ'. The 'AC' column (SN) contains 'SN'. The 'BL' column (BL) contains 'BL'. The spreadsheet also shows the Microsoft Office ribbon with various tabs like 'Soubor', 'Domů', 'Vložení', 'Rozložení stránky', 'Vzorce', 'Data', 'Revize', and 'Zobrazení'. The status bar at the bottom shows the page number '1' and the date '22/09/2016'.

**Proteinové sekvence????**



## Proteinové sekvence

Proteinové sekvence jsou odvozeny přepisem kódujících DNA sekvencí (CDS)

Každý region DNA má 6 čtecích rámců

5' 3'  
atgcccaagctgaatagcgtagaggggtttcatcatttgaggacgatgtataa

Čtecí rámeček začíná **atg (Met)** u většiny druhů a končí stop kodonem (**taa, tag or tga**).

# Proteinové sekvence

Proteinové sekvence jsou odvozeny přepisem kódujících DNA sekvencí (CDS)

Každý region DNA má 6 čtecích rámců, 3 v každém směru

5' 3'  
**atg**ccaagctgaatagcgtagaggggtttcatcattgaggacgatgtataa

**1** **atg** ccc aag ctg aat agc gta gag ggg ttt tca tca ttt gag gac gat gta **taa**  
M P K L N S V E G F S S F E D D V \*

**2** **tgc** cca agc **tga** ata gcg tag agg ggt ttt cat cat ttg agg acg atg tat  
C P S \* I A \* R G F H H L R T M Y

**3** **gcc** caa gct gaa **tag** cgt aga ggg gtt ttc atc att **tga** gga cga tgt ata  
A Q A E \* R R G V F I I \* G R C I

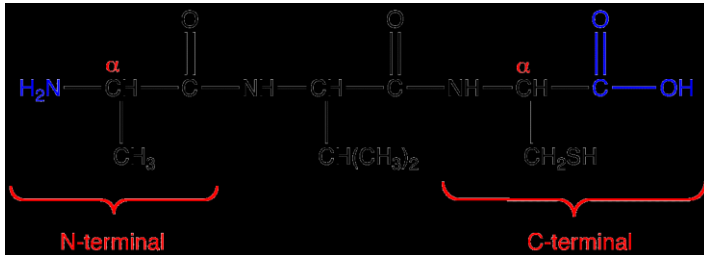
Čtecí rámec začíná **atg (Met)** u většiny druhů a končí stop kodonem (**taa, tag or tga**).



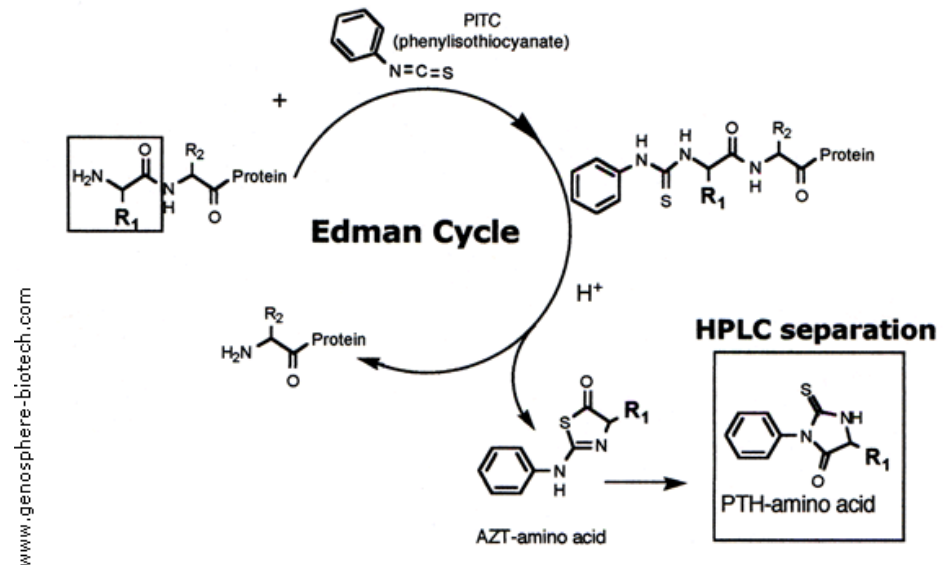
# Proteinové sekvence

Další zdroje proteinových sekvencí

- přímé sekvenování Edmanova degradace (odbourávání aminokyselin z N-konce a jejich identifikace)



## N-terminal sequencing cycle:



+ MS/MS

„nepřímé“ sekvenování MS/MS experimenty  
Získaná spektra se porovnávají vůči  
databázi

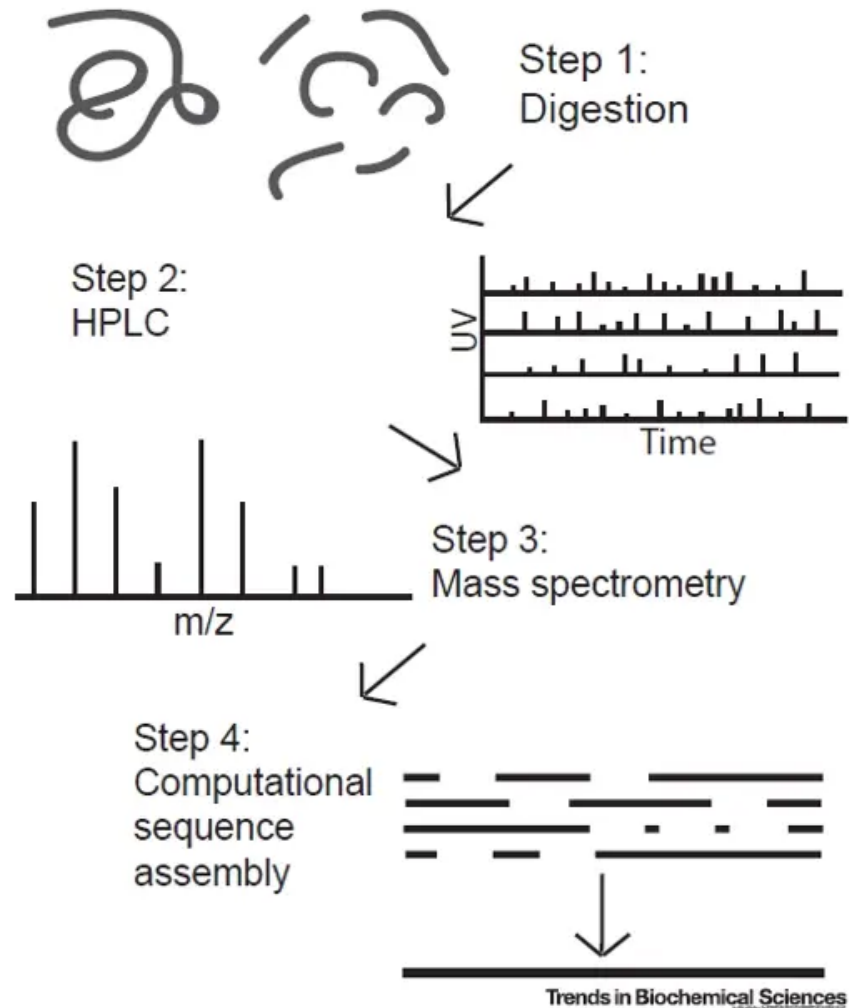
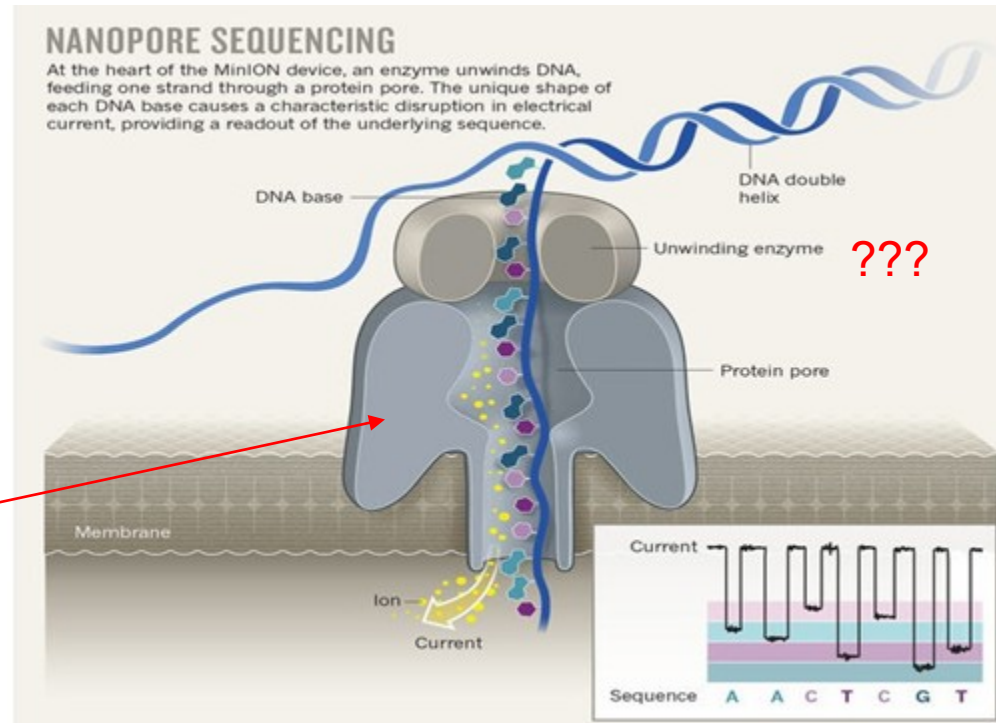
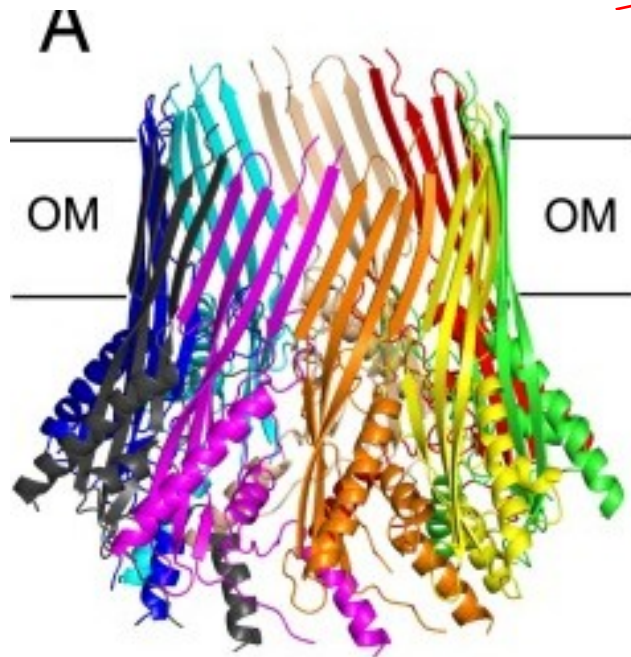


Figure 1. Typical Mass Spectrometry Workflow [1]

# DNA Nanopore sequencing

pore-forming proteins, such as an engineered version of CsgG



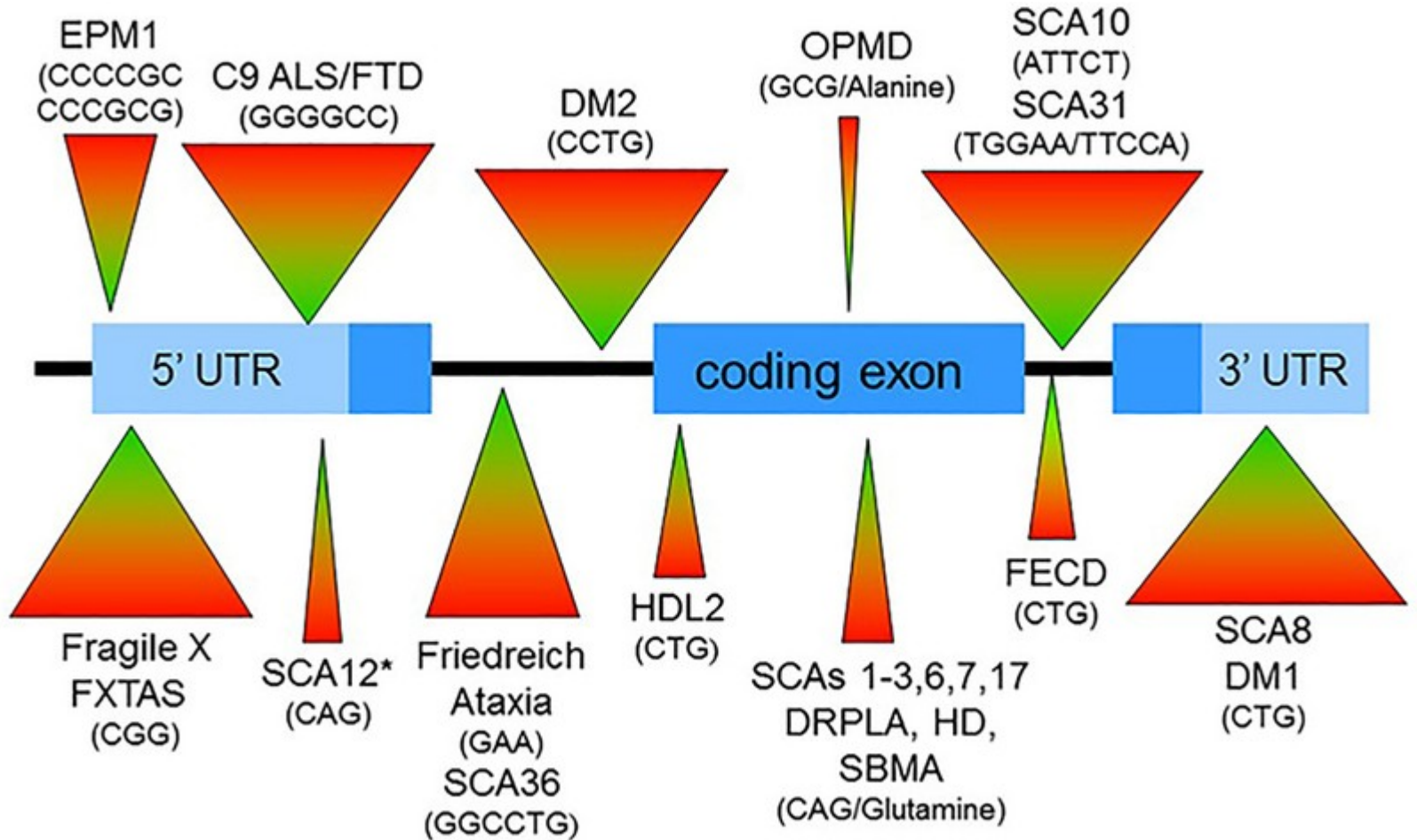
[https://www.youtube.com/watch?v=RcP85JHLmnl&ab\\_channel=OxfordNanoporeTechnologies](https://www.youtube.com/watch?v=RcP85JHLmnl&ab_channel=OxfordNanoporeTechnologies)

## Detekce expanzi pomocí nanoporového sekvenování

Repeat expansions cause many neurologic diseases

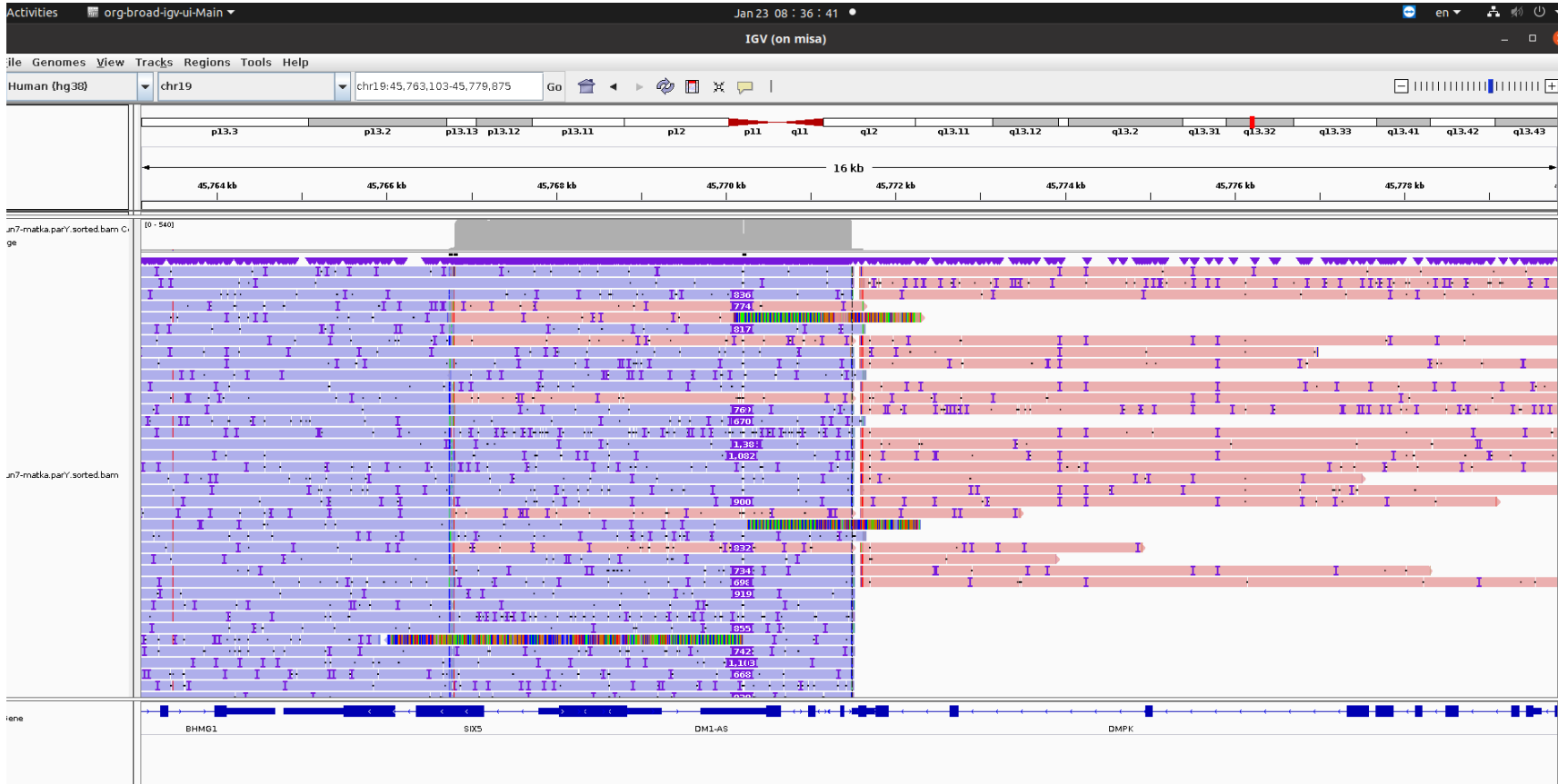
- CAG – at least 10 diseases (Huntington disease, spinal and bulbar muscular atrophy, dentatorubral pallidoluysian atrophy and seven SCAs)
- CGG – fragile X, fragile X tremor ataxia syndrome, other fragile sites (GCC, CCG)
- CTG – myotonic dystrophy type 1, Huntington disease-like 2, spinocerebellar ataxia type 8, Fuchs corneal dystrophy
- GAA – Friedreich ataxia
- GCC – FRAXE mental retardation
- GCG – oculopharyngeal muscular dystrophy
- CCTG – myotonic dystrophy type 1
- ATTCT – spinocerebellar ataxia type 10
- TGGAA – spinocerebellar ataxia type 31
- GGCCTG – spinocerebellar ataxia type 36
- GGGGCC – C9ORF72 frontotemporal dementia/amyotrophic lateral sclerosis
- CCCC GCCCGCG – EPM1 (myoclonic epilepsy)

# Lokalizace expanzí v genu

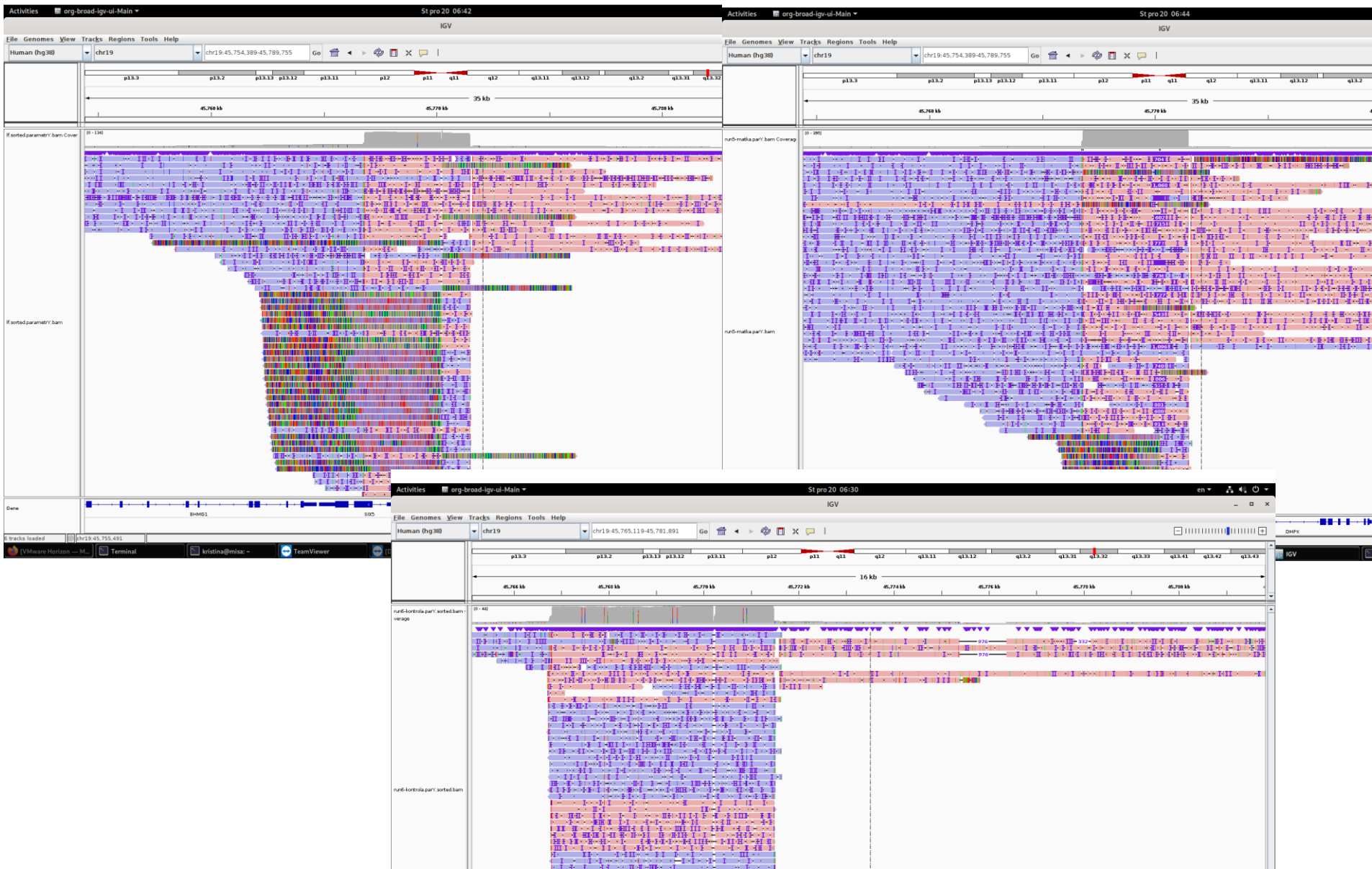


markedly different sizes of pathogenic expansions are suggested by the varying sized triangles

# Expanze repetit CTG v genu DMPK – myotonicka dystrofie

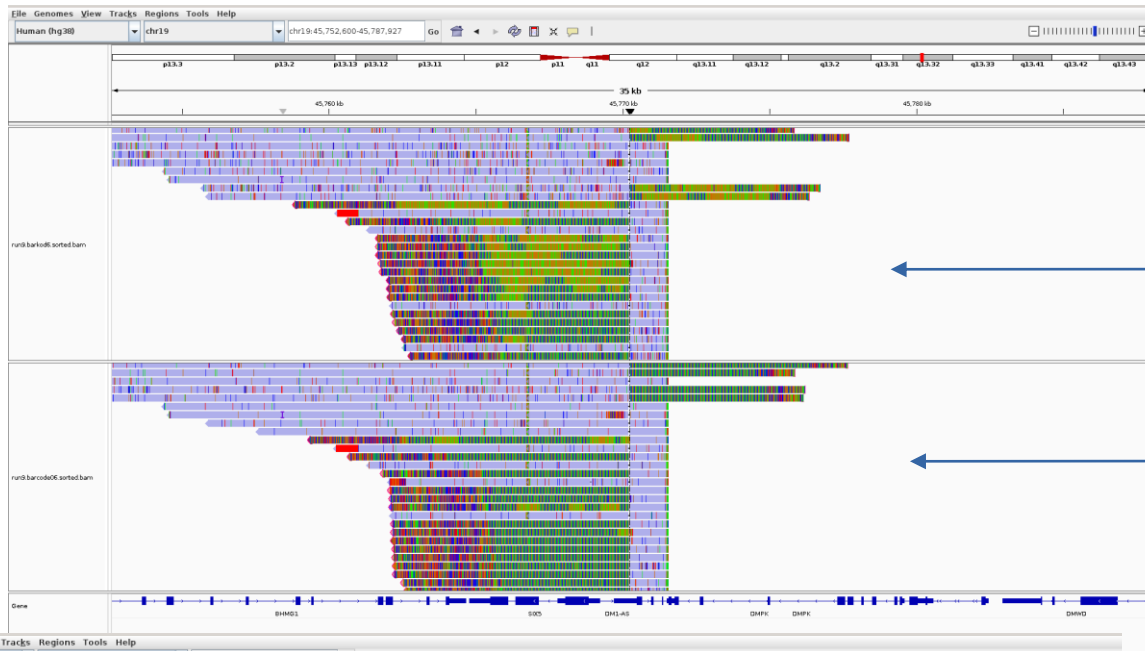


# Expanze repetit CTG v genu DMPK – myotonická dystrofie



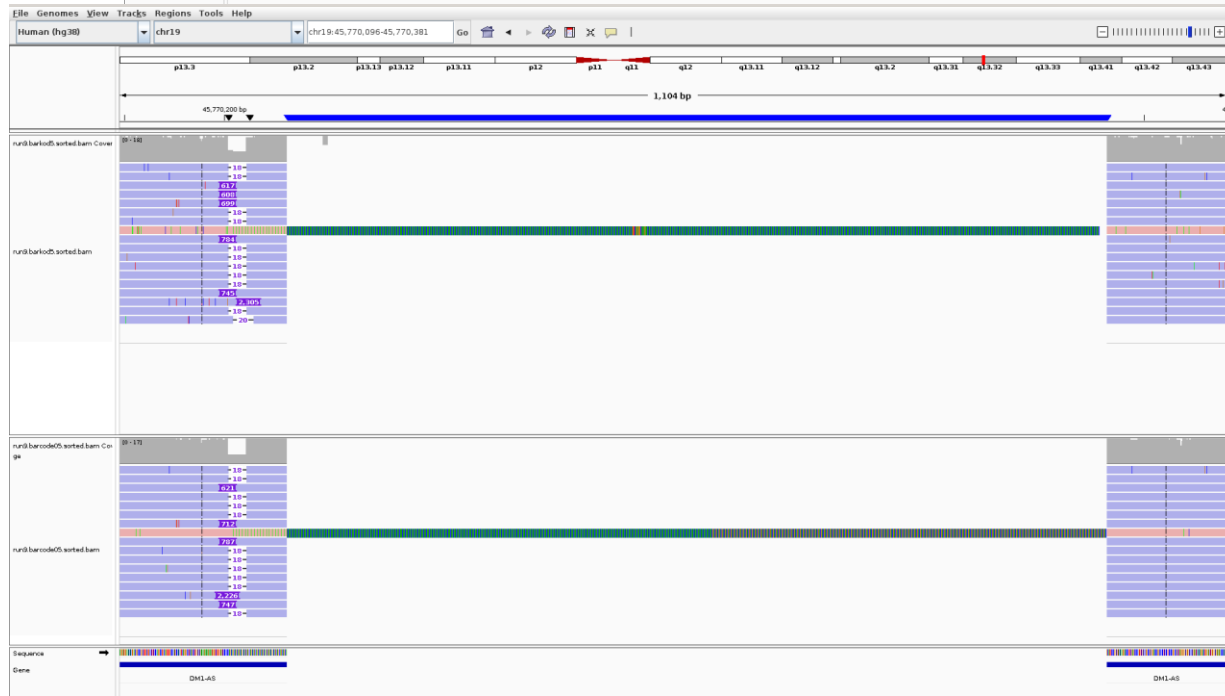


# Problém při basecallingu repetitivních oblastí



Základní  
basecalling

Zpřesněný basecalling

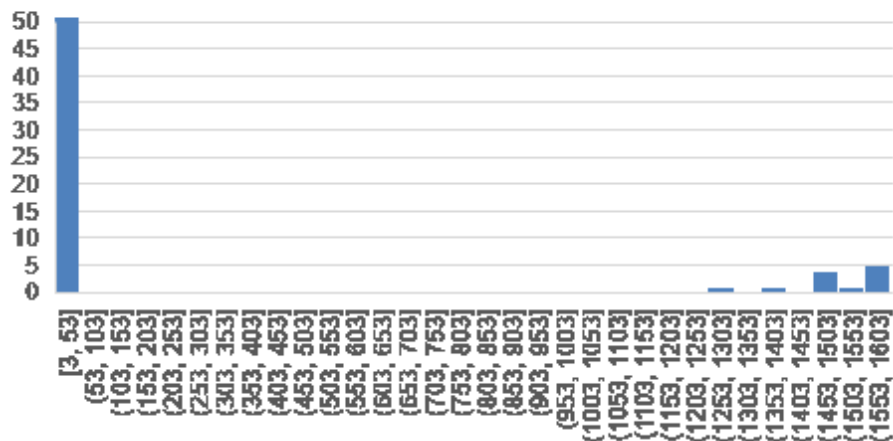


Základní  
basecalling

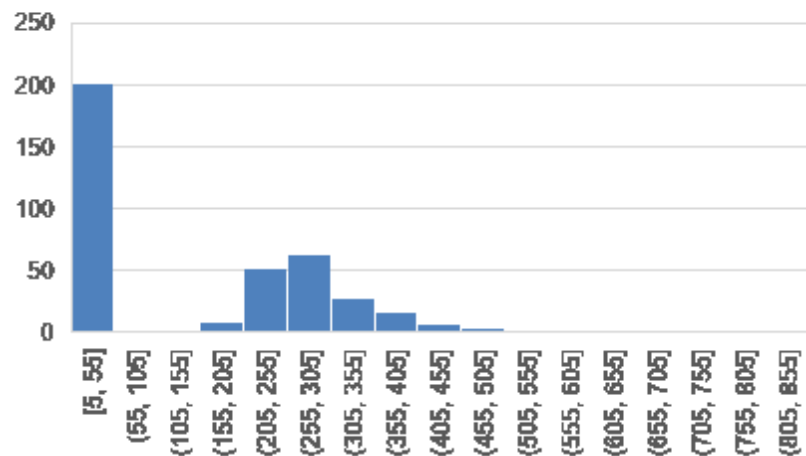
Zpřesněný basecalling



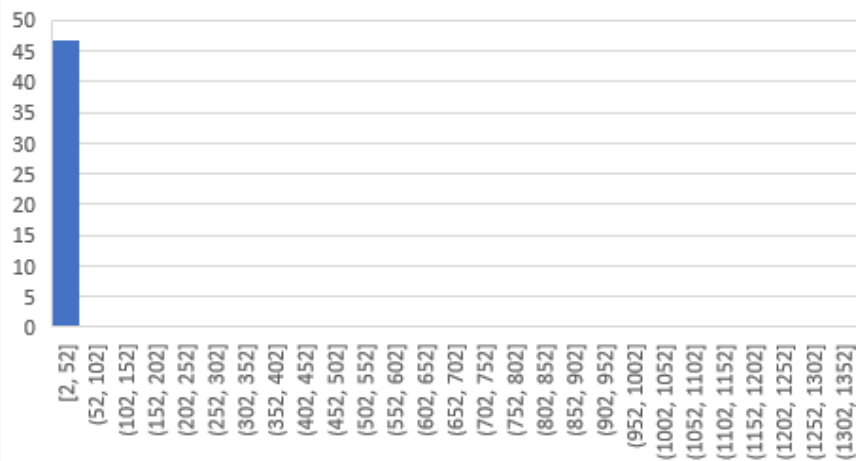
### proband - četnost CAG v DMPK



### matka - četnost CAG v DMPK



### kontrola - četnost CAG v DMPK



# Protein Nanopore sequencing???

The screenshot shows a web browser window displaying a Nature Biotechnology article. The browser's address bar shows the URL: <https://www.nature.com/articles/s41587-019-0401-y?proof=t>. The article title is "Reading amino acids in a nanopore" by Stefan Howorka & Zuzanna S. Siwy, published on 23 January 2020. The article is categorized under "SEQUENCING TECHNOLOGY". The abstract reads: "In a step toward nanopore sequencing of proteins, an aerolysin pore discriminates many of the proteinogenic amino acids." Below the abstract, it states: "Nanopore sequencing of DNA is performed by threading single-stranded DNA (ssDNA) through a narrow pore and measuring electrical signatures of the four bases – tiny". On the right side of the page, there is a "Download PDF" button and a section for "Associated Content" featuring a letter titled "Electrical recognition of the twenty proteinogenic amino acids using an aerolysin nanopore" by Hadjer Ouldali, Kumar Sarthak, and Abdelghani Oukhaled. The browser's taskbar at the bottom shows various application icons and the system clock indicating 10:19 AM on 4/15/2021.

## Main issues ??

Unlike 4 nucleotides that make up DNA, proteins are composed of 20 amino acids, each with a different charge and virtually 100s of potential post-translational modifications [2]. Moreover, it has proved difficult to control the translocation of peptides across a nanopore as enzymes capable of ratcheting the peptide at a controlled rate have been difficult to identify and conjugate with existing nanopores.

- Mutace/varianty (odchylky od referenční sekvence)

NGS:

panel genů

Exom → WES – whole exom sequencing (2% genomu)

Genom → WGS – whole genome sequencing

Sekvenování identifikuje (známé) varianty nebo odhalí nové

(jsou skutečně kauzální???)

Je nalezená mutace  
popsaná????



mutační databáze: germinální mutací (**HGMD**), somatické (COSMIC)

Specializované: (**PAHdb**, **Leiden**, **IARC TP53 Database – somatické/germinální**)

## Human Gene Mutation Database (HGMD) –

databáze mutací  
lidských genů v  
kódujících oblastech.  
(vychází z primárních  
dat, manuální a  
Automatické  
prohledávání)

my.qiagen.com/bbp/view/hgmd/pro/stats.php

FN Brno Mail - Rěblová Kam... Historie QIAGEN Digital Insi... Seznam benefitů F... Sign out Mail - Rěblová Kam... Pošta - Rěblová Ka...

### HGMD® Professional 2023.3

Gene Mutation Phenotype Reference Batch Advanced | Statistics Information Support | Home Logout

#### HGMD Statistics

Data type	Total number of entries for release 2023.3	Total number of entries for release 2023.2
<a href="#">HGMD genes</a> <small>(plus alternative isoforms)</small>	16539 (1134)	16413
<a href="#">HGMD cDNA sequences</a>	17481	17312
<a href="#">HGMD primary references</a>	91058	89689
<a href="#">HGMD additional references</a>	64926	63664
<a href="#">HGMD phenotypes</a>	35868	34969
<a href="#">HGMD mapped phenotypes</a>	36056	35143

#### Number of entries by mutation type

Mutation type	Total number of entries for release 2023.3 (disease-associated functional polymorphism sub-total)	Total number of entries for release 2023.2
<a href="#">Missense/nonsense</a>	276726 (8521)	268697
<a href="#">Splicing</a>	39662 (846)	38443
<a href="#">Regulatory</a>	6412 (3334)	6225
<a href="#">Small deletions</a>	62838 (445)	61216
<a href="#">Small insertions</a>	27242 (249)	26480
<a href="#">Small indels</a>	5231 (79)	5101
<a href="#">Repeat variations</a>	692 (370)	687
<a href="#">Gross insertions/duplications</a>	7148 (284)	7043
<a href="#">Complex rearrangements</a>	2869 (146)	2818
<a href="#">Gross deletions</a>	27882 (194)	27375
<i>Total</i>	<i>456702 (14452)</i>	<i>444085</i>

#### Number of entries by variant class

Variant class	Total number of entries for release 2023.3	Total number of entries for release 2023.2
<a href="#">SNP</a>	297490	291395
<a href="#">indel</a>	144237	137768
<a href="#">SV</a>	5584	5555

Co když mutace není popsána?

Jak rozlišit mezi:

Bežnou variantou [**single nucleotide polymorphisms (SNPs)**]

**a**

**kauzální mutací?**

## Přístupy studia efektů missense mutací:

- Funkční analýzy (exprese proteinu, stabilita, specifické charakteristiky)  
časově a finančně náročné
- Automatické *in silico* programy (SIFT, PolyPhen, SNPs3D, FODLX..)  
rychlé a spolehlivost průměrně 70 %
- Strukturní analýza pomocí molekulového modelování  
časově a finančně středně náročné, pohled do mechanismu účinky na  
atomární úrovni, je to ale stále predikce

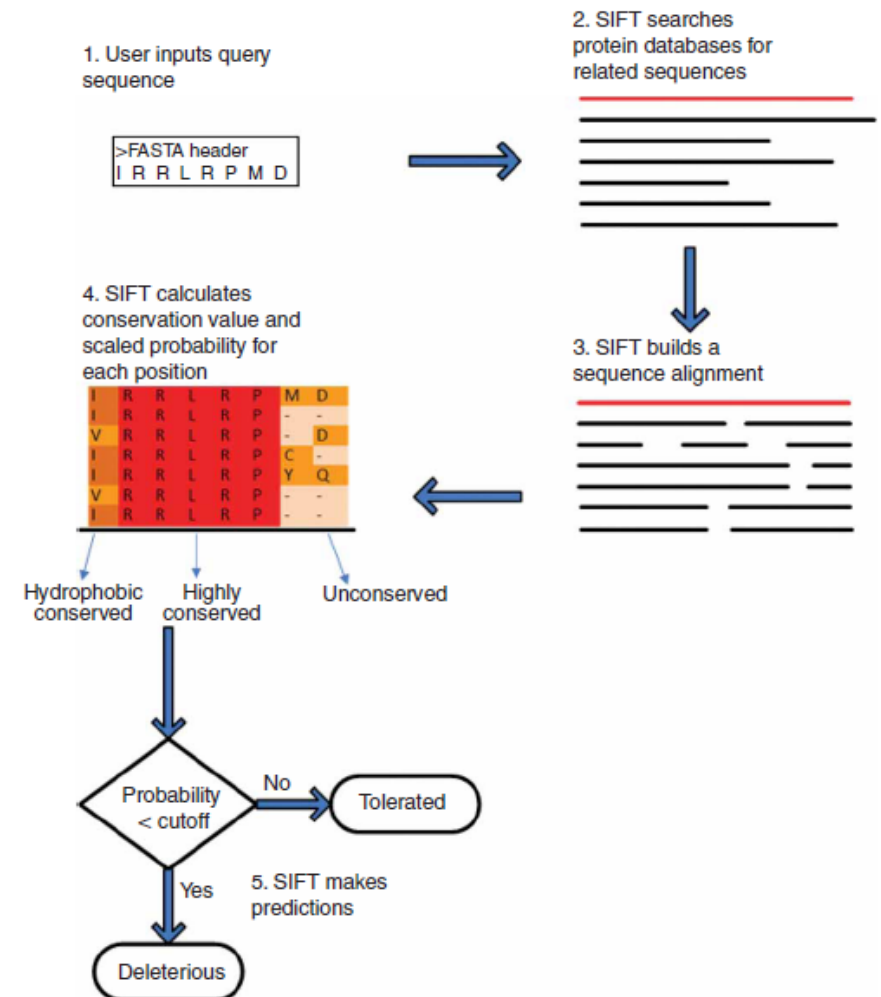
# Automatické analýzy missense mutací pomocí bioinformatických nástrojů

**SIFT** - Sorting Tolerant From Intolerant  
(<http://sift.jcvi.org>)

SIFT assumes that important positions in a protein sequence have been conserved throughout evolution and therefore substitutions at these positions may affect protein function.

Using sequence homology, SIFT predicts the effects of all possible substitutions at each position in the protein sequence.

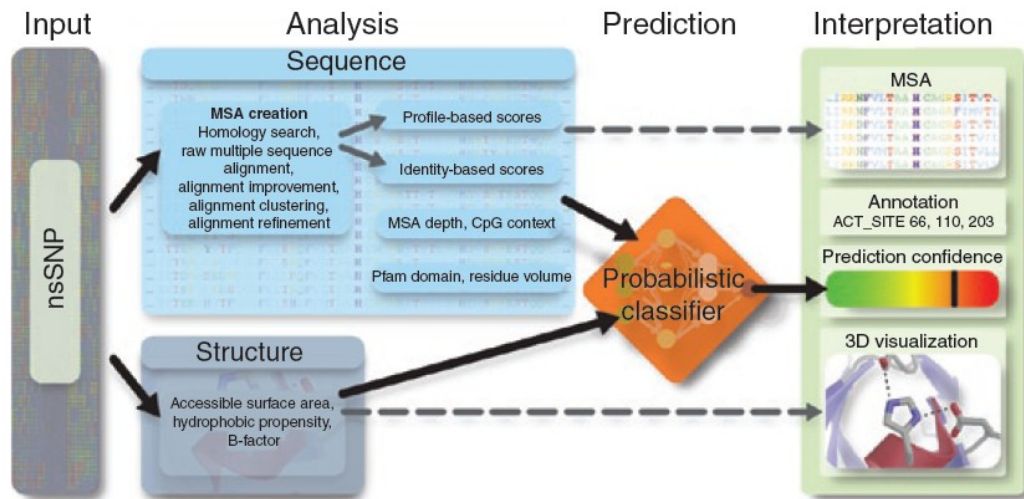
Úspěšnost predikce pro skupinu kauzálních mutací 69%.  
Falešně pozitivní 19 %.



# Automatické analýzy missense mutací pomocí bioinformatických nástrojů

## Polyphen -2

PolyPhen-2 uses eight sequence-based and three structure-based predictive features



Structural features. Three additional features were selected for proteins with known 3D structures: 1) the accessible surface area of the wild-type amino acid residue, 2) the change in the hydrophobic propensity in the form of “knowledge-based potential”, and 3) crystallographic B-factor reflecting conformational mobility of the wild-type amino acid residue15.

Úspěšnost predikce 73%

Falešně pozitivní ~ 20 %.



## Automatické analýzy missense mutací pomocí bioinformatických nástrojů

**FOLDX** (<http://foldx.embl.de/>)

Změna celkové volné energie

$$\Delta\Delta G = \Delta G_{mut} - \Delta G_{wt}$$

$\Delta G_{mut}$  - free energy difference between the folded and unfolded states

Empirical force field that was developed for the rapid evaluation of the effect of mutations on the stability, folding and dynamics of proteins and nucleic acids

$$\begin{aligned}\Delta G = & a \cdot \Delta G_{vdw} + b \cdot \Delta G_{solvH} + c \cdot \Delta G_{solvP} + d \cdot \Delta G_{wb} \\ & + e \cdot \Delta G_{hbond} + f \cdot \Delta G_{el} + g \cdot \Delta G_{kon} + h \cdot T\Delta S_{mc} \\ & + k \cdot T\Delta S_{sc} + l \cdot \Delta G_{clash}.\end{aligned}$$

(a . . . l) are relative weights of the different energy terms used for the free energy calculation

Výpočty na lokálním PC ne přes web rozhraní nyní.

Lze měnit počáteční nastavení: T(K), c(M)...

## Automatické analýzy missense mutací pomocí bioinformatických nástrojů

**FOLDX** (<http://foldx.embl.de/>)

<i>FoldX component</i>	<i>Energy (kcal.mol-1)</i>
Backbone Hbond	-48.30
Sidechain Hbond	-18.40
Van der Waals	-62.33
Electrostatics	-2.79
Solvation Polar	84.34
Solvation Hydrophobic	-76.71
Van de Waals clashes	0.00
entropy side chain	34.93
entropy main chain	81.94
sloop_entropy	0.00
mloop_entropy	0.00
cis_bond	0.00
torsional clash	0.40
backbone clash	1.20
helix dipole	0.00
water bridge	0.00
disulfide 0.00	
electrostatic kon	0.00
partial covalent bonds	0.00
Total	-5.72

$$\Delta\Delta G(\text{change}) = \Delta G(\text{MT}) - \Delta G(\text{WT})$$

$\Delta\Delta G(\text{change}) > 1\text{kcal/mol}$ : the mutation is destabilizing

$\Delta\Delta G(\text{change}) < -1\text{ kcal/mol}$ : the mutation is stabilizing

Úspěšnost predikce 60%

# Konsensus mutačních programů..... např. Varsome

The screenshot displays the Varsome website interface for a specific variant. The browser address bar shows the URL: <https://varsome.com/variant/hg19/chr19-11211462-G-A?>. The page title is "chr19-11211462-G-A". The navigation bar includes "varsome", "Search", "Editions", "About", "Community", "News", "Demo", "Sign in", and "Join".

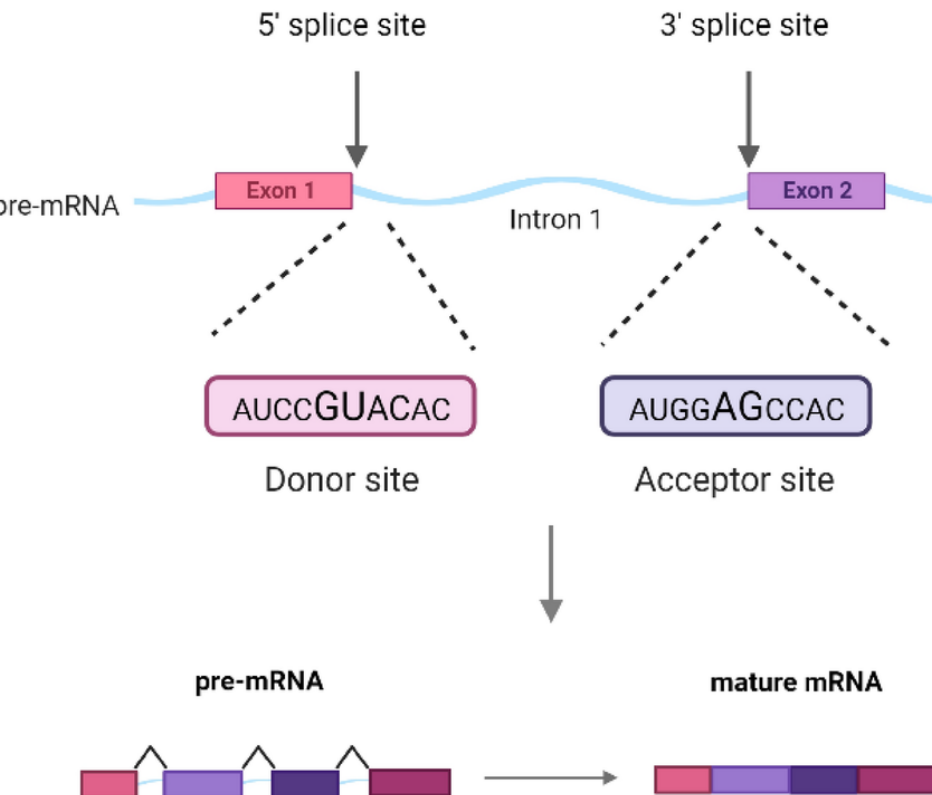
The main content area is a grid of data panels for the variant "chr19-11211462-G-A". The panels include:

- General Information: SNV, LDLR(NM\_000527.5):c.190+441G>A
- PharmGKB: No data available
- Germline Classification: Labeled "Likely Benign" with a score of -4 points (0 P - 4 B).
- Frequencies: exomes: not found, genomes:  $f = 0$  (cov: 19.6 low)
- Conservation Scores: phyloP100: -1.031
- Structural Variants: No data available
- Genes: LDLR
- Transcripts: NM\_000527.5 - non coding, MANE Select
- ClinVar: No data available
- MitoMap: No data available
- In-Silico Predictors: BP4: Benign Strong (score 4)
- Beacon Network: No data available
- Community Contributions: Region Browser
- LOVD: No data available
- Deafness Variation Database: No data available
- ClinGen: No data available
- Protein Viewer: No data available
- Publications: Variant: 0, Gene: 2098
- Expression Data: No data available
- Uniprot Variants: No data available
- OMIM: No data available
- GWAS: No data available

Below the grid is the "Variant" section, which provides a detailed view of the variant's genomic context. It includes fields for Chromosome (chr19), Position (11211462), REF Sequence (G), ALT Sequence (A), Variant type (SNV), Cytoband (19p13.2), HGVS (LDLR(NM\_000527.5):c.190+441G>A), RS ID (rs1288832894), and Gene symbol (LDLR). A button "Connect with past and future viewers of this variant..." is visible. At the bottom, a note states: "VarSome.com is for research use only. Find out about our clinically certified platform: VarSome Clinical."

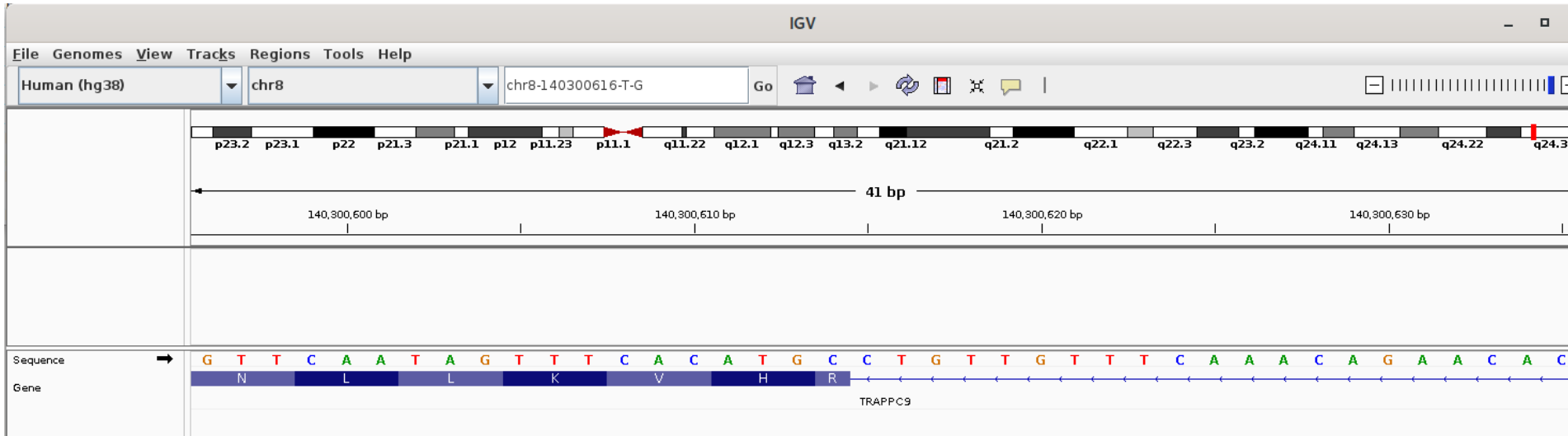
MetaRNN and MetaRNN-indel are pathogenicity prediction scores for human nonsynonymous SNVs (nsSNVs) and non-frameshift (NF) indels. They integrated information from 28 high-level annotation scores (16 functional prediction scores including SIFT, Polyphen2\_HDIV, Polyphen2\_HVAR, MutationAssessor, PROVEAN, VEST4, M-CAP, REVEL, MutPred, MVP, PrimateAI, DEOGEN2, CADD, fathmm-XF, Eigen and GenoCanyon, 8 conservation scores including GERP, phyloP100way Vertebrate, phyloP30way\_mammalian, phyloP17way\_primate, phastCons100way Vertebrate, phastCons30way\_mammalian, phastCons17way\_primate and SiPhy, and 4 allele frequency information from the 1000 Genomes Project, ExAC, gnomAD exome, and gnomAD genome) and produce an ensemble prediction model with a deep recurrent neural network (RNN). The final prediction is the likelihood of a nsSNV or NF indel being pathogenic.

# Predikční nástroje na sestřihové varianty - nejčastěji rozhraní exon/intron



Small nuclear RNAs (snRNAs) are critical components of the spliceosome that catalyze the splicing of pre-mRNA.

# Predikční nástroje na rozhraní exon/intron – sestřihové varianty



Examples (on hg38):

chr8-140300616-T-G  
 6 31740453 G T  
 NM\_001089.3(ABCA3):c.875A>T (p.Glu292Val)

## SPLICE AI

[\[show more examples\]](#)

chr8-140300616-T-G

Genome version:  hg19  hg38 Gencode:  basic  comprehensive ?

Max distance:  ?  masked scores ?  REF & ALT scores ?

Submit

November 3, 2024

- added choice of Gencode [basic](#) or [comprehensive](#) transcripts  
 - updated to Gencode v47

June 20, 2024

- fixed SpliceAI visualizations to clarify positions of predicted splicing changes. See [issue #70](#) for details.

[\[show older updates\]](#)

Related web tools:

[liftover](#): for variants/positions/intervals (hg19 <=> hg38 <=> T2T)  
[TGG Viewer](#): igv.js-based web viewer for public reference tracks and private data in Google Storage buckets. Has custom track types for RNA-seq [splice junctions](#) and [gCNV](#) variants.

SpliceAI scores: ?

Variant	Gene	<input type="checkbox"/> = MANE Select transcript <input type="checkbox"/> = non-coding transcript	Δ type	Δ score ?	position ?
chr8-140300616-T-G splice acceptor variant UCSC, gnomAD	TRAPP9 (ENSG00000167632.18 / ENST00000438773.4 / NM_001160372.4) protein coding MANE Select transcript (minus strand) OMIM, GTEx, gnomAD, ClinGen, Ensembl, Decipher, GeneCards	<input type="checkbox"/>	Acceptor Loss	0.83	-2 bp
		<input type="checkbox"/>	Donor Loss	0.62	-147 bp
		<input type="checkbox"/>	Acceptor Gain	0.04	-32 bp
		<input type="checkbox"/>	Donor Gain	0.00	

# Strukturní analýza pomocí molekulového modelování

Nutná podmínka – xray nebo nmr struktura!

Visualizační program: VMD, PYMOL



The screenshot shows the homepage of the VMD website. At the top, it identifies the 'THEORETICAL and COMPUTATIONAL BIOPHYSICS GROUP' at the 'UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN'. The main content area features the VMD logo and a description: 'VMD is a molecular visualization program for displaying, animating, and analyzing large biomolecular systems using 3-D graphics and built-in scripting.' A 'Spotlight' section highlights a memorial issue in the Journal of Physical Chemistry from April 20th, 2017, honoring Klaus Schulten. A sidebar on the left contains navigation links for Home, Research, Publications, Software, Instruction, News, Galleries, Facilities, and About Us. A 'VMD Mailing' link is also present at the bottom left.

The screenshot shows the PyMOL website homepage. The header includes 'PyMOL by Schrödinger' and navigation links for 'DOWNLOAD', 'SCREENSHOTS', 'SUPPORT', and 'CONTACT'. The main visual is a large, detailed 3D molecular model of a protein structure, rendered in a dark blue and purple color scheme. The text 'Introducing PyMOL 3.1' is prominently displayed. Below this, a paragraph states: 'PyMOL is a user-sponsored molecular visualization system on an open-source foundation, maintained and distributed by Schrödinger.' At the bottom, there are three buttons: 'DOWNLOAD NOW', 'BUY LICENSE', and 'PYMOL 3 PRODUCT PAGE'.

# 3D struktury proteinů – databáze PDB

The image shows a screenshot of the RCSB PDB website homepage. The browser window has a single tab titled "RCSB PDB: Homepage" and the address bar shows "https://www.rcsb.org". The website header includes navigation menus for "RCSB PDB", "Deposit", "Search", "Visualize", "Analyze", "Download", "Learn", "About", "Documentation", "Careers", and "COVID-19". There are also buttons for "MyPDB" and "Contact us".

The main content area features the RCSB PDB logo and statistics: "212,599 Structures from the PDB" and "1,068,577 Computed Structure Models (CSM)". A search bar is present with the placeholder text "Enter search term(s), Entry ID(s), or sequence" and an "Include CSM" toggle switch. Below the search bar are links for "Advanced Search" and "Browse Annotations", and a "Help" link.

A banner below the search bar reads "New: More Computed Structure Models (CSM) available" with a "Learn more" button. The left sidebar contains navigation options: "Welcome", "Deposit", "Search", "Visualize", and "Analyze".

The main content area contains the following text:

RCSB Protein Data Bank (RCSB PDB) enables breakthroughs in science and education by providing access and tools for exploration, visualization, and analysis of:

- Experimentally-determined 3D structures from the **Protein Data Bank (PDB)** archive
- Computed Structure Models (CSM)** from AlphaFold DB and ModelArchive

These data can be explored in context of external annotations providing a structural view of biology.

On the right side, there is a section titled "November Molecule of the Month" featuring a 3D molecular structure visualization of a protein complex, colored in shades of blue, green, and yellow.

The Windows taskbar at the bottom shows the search bar with "Type here to search", several application icons, the system tray with a temperature of 6°C, and the date and time: 10:26 AM 11/24/2023.





Mutace narušuje  
strukturu proteinu

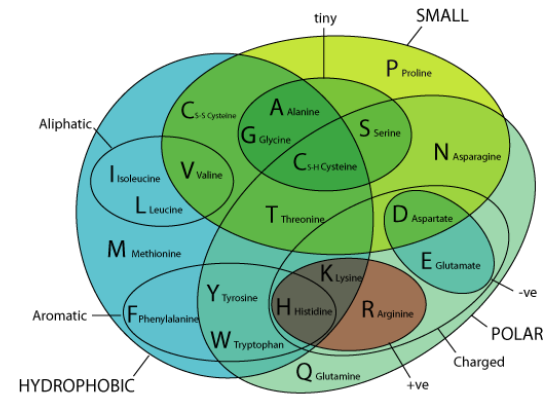


Vytvářejí síť kontaktů s okolními aminokyselinami  
*H-bonding, stacking, salt bridges*



narušení struktury proteinu vede k destabilizaci proteinu,  
unfoldingu příp. agregaci

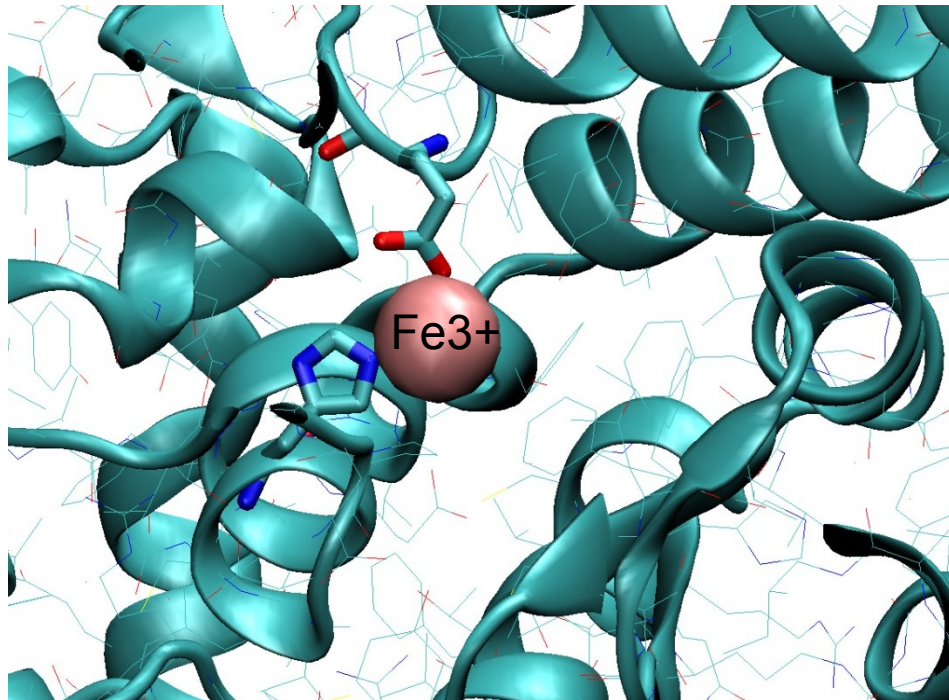
*změna náboje,  
polarity  
velikost aminokyseliny*



Mutace ve funkčních  
místech proteinu



(např. vazbu ligandu/kofaktoru, příp. vazba specifických iontů součást konformace aktivního místa)



Obr. PAH

## Ukazatele kauzality

### Hodnocené znaky:

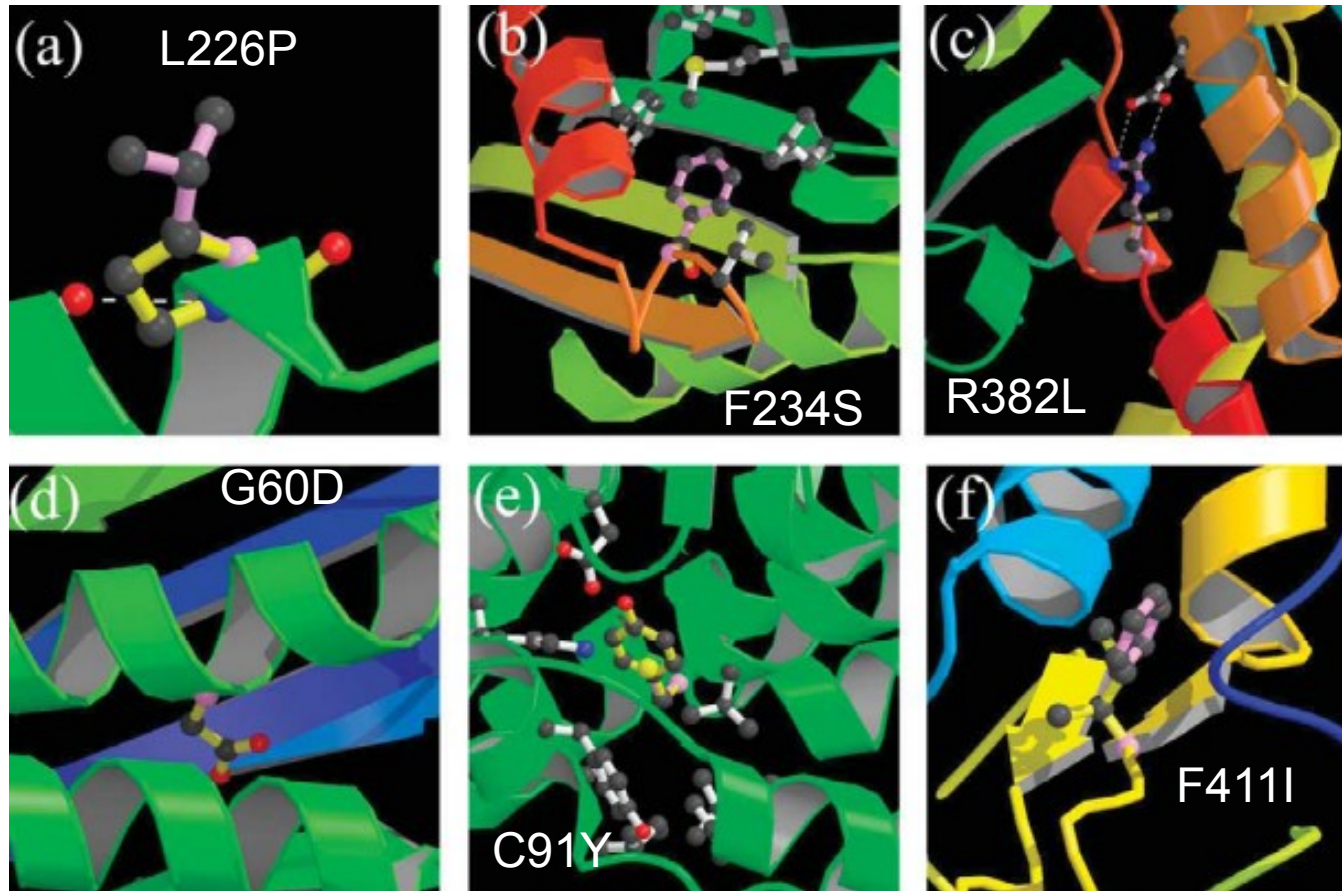
- Specifické kontakty postraními řetězci AA
- Výskyt AA v aktivním místě
- Zanořenost AA v proteinu
- změna objemu AA
- změna náboje
- změna polarity
- Konzervovanost AA
- Přítomnost helix/turn breakers

1. Kauzální mutace uvnitř proteinu, polymorfismy spíše na povrchu

2. Kauzální mutace: wt mutace častěji vytváří vodíkové vazby a jiné kontakty



Examples of disease caused by structure destabilizing factors:  
In silico analysis of mutations



Causal mutations:

bonds of wild-type side-chains are shown purple,  
and bonds of the mutant side-chains are yellow.