

MUNI
SCI



MASARYK UNIVERSITY
FACULTY OF SCIENCE
DEPARTMENT OF EXPERIMENTAL BIOLOGY
LABORATORY OF MICROBIAL MOLECULAR DIAGNOSTICS

Bi5000 – Bioinformatika
Lekce 02

Manipulace se sekvenčními daty

Sylabus výuky předmětů Bi5000 Bioinformatika a Bi5000c Bioinformatika - cvičení v semestru podzim 2024



Přednášky:	úterý 11:00-12:50 prezenčně v učebně B11-306
Kontakt:	prof. Mgr. Jiří Damborský, Dr. (1441@mail.muni.cz) prof. RNDr. Roman Pantůček, Ph.D. (pantucek@sci.muni.cz)
Cvičení:	středa 10:00-17:00 prezenčně v učebně B09-316, dle seminářních skupin
Kontakt:	Ing. Miloš Musil (imusilm@fit.vutbr.cz) prof. RNDr. Roman Pantůček, Ph.D.

Datum	č.	Předmět	Lekce
17.09.2024			Výuka se nekoná z důvodu povodní
18.09.2024			
24.09.2024	2	Bi5000	Manipulace se sekvenčními daty
25.09.2024		Bi5000c	Cvičení - textové vyhledávání v databázích
01.10.2024	1	Bi5000	Bioinformatika – základní definice, molekulárně biologické databáze
02.10.2024		Bi5000c	Cvičení – formáty sekvencí, manipulace se sekvenčními daty
08.10.2024	3	Bi5000	Posuzování podobnosti sekvencí nukleových kyselin a proteinů
09.10.2024		Bi5000c	Cvičení – párové přiložení sekvencí, BLAST
15.10.2024	4	Bi5000	Mnohonásobné přiložení sekvencí a fylogeneze
16.10.2024		Bi5000c	Cvičení - mnohonásobné přiložení a fylogeneze
22.10.2024	5	Bi5000	Genomové projekty, sekvenování nové generace
23.10.2024		Bi5000c	Cvičení – návrh oligonukleotidů pro PCR, sekvenování, klonování a mutagenizi
29.10.2024	6	Bi5000	Počítačové vyhledávání genů a srovnávací genomika
30.10.2024		Bi5000c	Cvičení – hledání prokaryotických genů, práce s programy pro srovnávací genomiku
05.11.2024	7	Bi5000	Analýza sekvencí proteinů
06.11.2024		Bi5000c	Cvičení – analýza sekvencí proteinů
12.11.2024	8	Bi5000	Strukturní databáze
13.11.2024		Bi5000c	Cvičení – základní strukturní analýzy, vizualizace proteinových struktur v PyMOL
19.11.2024	9	Bi5000	Predikce struktury proteinů
20.11.2024		Bi5000c	Cvičení – predikce struktury proteinů
26.11.2024	10	Bi5000	Příprava anotované sekvence DNA pro zaslání do databáze
27.11.2024		Bi5000c	Cvičení – NGS data a lokální anotace dat
03.12.2024		Bi5000	1. Předtermín
04.12.2024		Bi5000c	Zápočtový test na PC (odpovědník v ISu)
10.12.2024		Bi5000	2. Předtermín, další termíny v lednu
11.12.2024		Bi5000c	Opravný zápočtový test na PC (odpovědník v ISu)

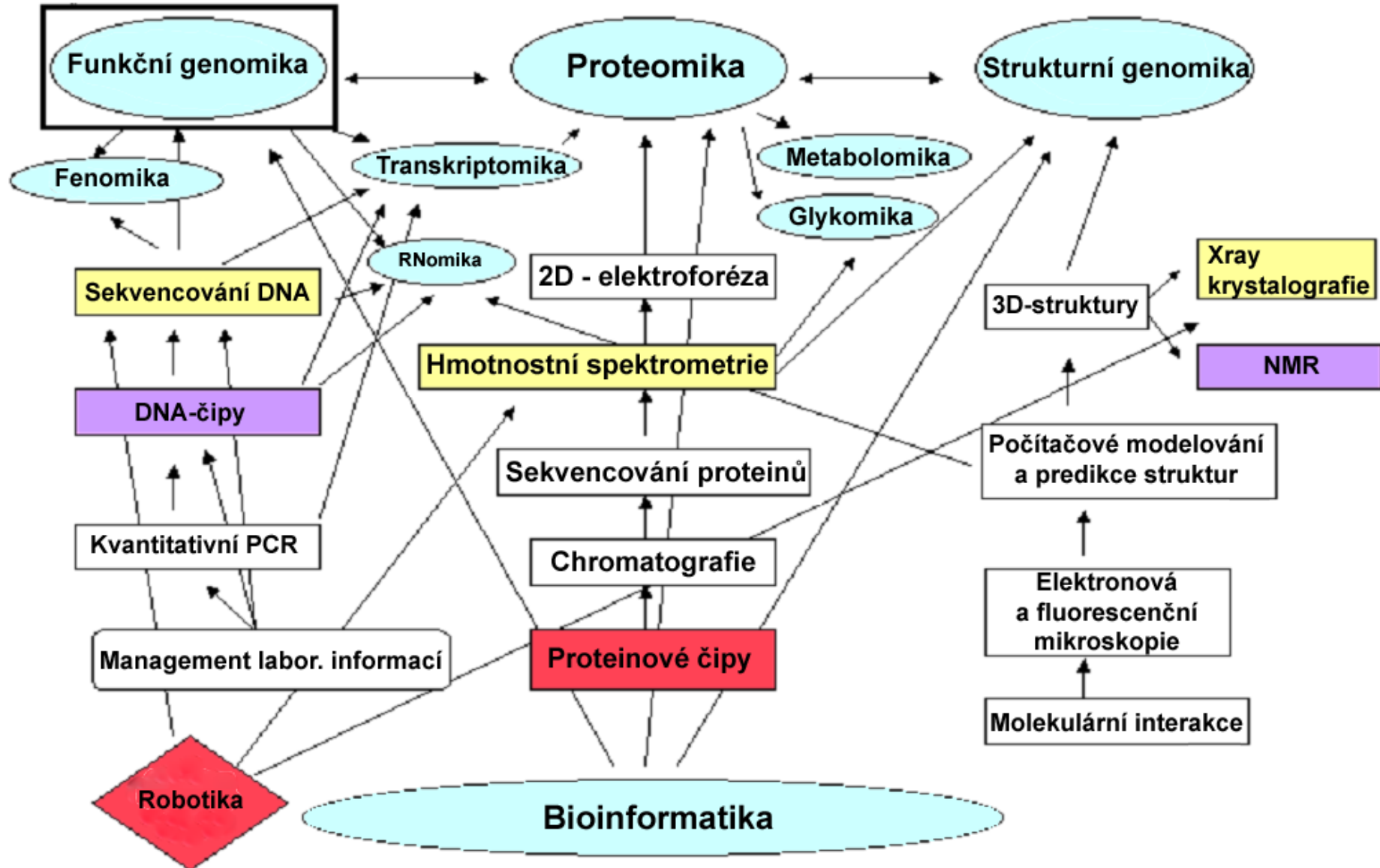
Bioinformatika je disciplína na rozhraní počítačových věd, informačních technologií, matematiky a biologie

- Termín bioinformatika se objevil poprvé v roce 1991
- Představuje spojení technologií z oblastí
molekulární biologie
informačních technologií
- Bioinformatika zahrnuje
studium a analýzu
praktické uchovávání
vyhledávání a zobrazování
modelování biologických dat
- Výpočetní nástroje umožňující analýzu dat a stanovení jejich vzájemných vztahů
- Dramatický nárůst množství dat a tím současně zvyšující se obtížnost jejich zkoumání a hodnocení ve vztahu k biologickým otázkám

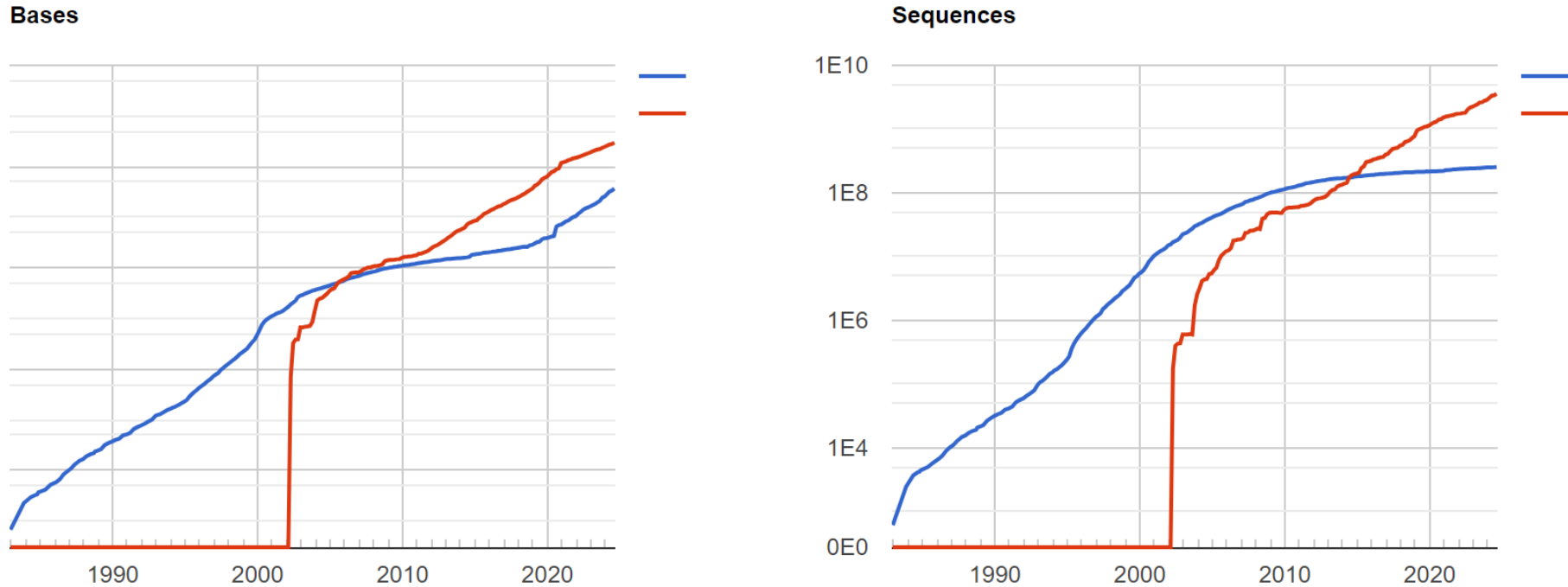
Základní zdroje a aplikace bioinformatiky

Výpočetní základy	Zdroje dat	Aplikace bioinformatiky
Algoritmy	Obecně dostupné databáze	Získávání dat
Grafika, vizualizace		Nástroje pro přístup k databázím
Zpracování signálu		Mapování a srovnávání genomů
Architektura hardwaru		Sekvenční příložen, assembly
Informační teorie		Identifikace genů
Správa databází		Funkční identifikace proteinů
Statistika		Molekulární evoluce
Simulace		Molekulární modelování
Umělá inteligence		Predikce struktur
Zpracování obrazu		Srovnávání struktur
Robotika	Zpracování laboratorních dat	Stanovení makromolekulárních struktur
Softwarové inženýrství		Vývoj léčiv na základě struktur

„-omiky“ v molekulární biologii



Trend nárůstu množství dat v bioinformatických databázích



Zdroj: <https://www.ncbi.nlm.nih.gov/genbank/statistics/>

Typy jednoduchých bioinformatických manipulací

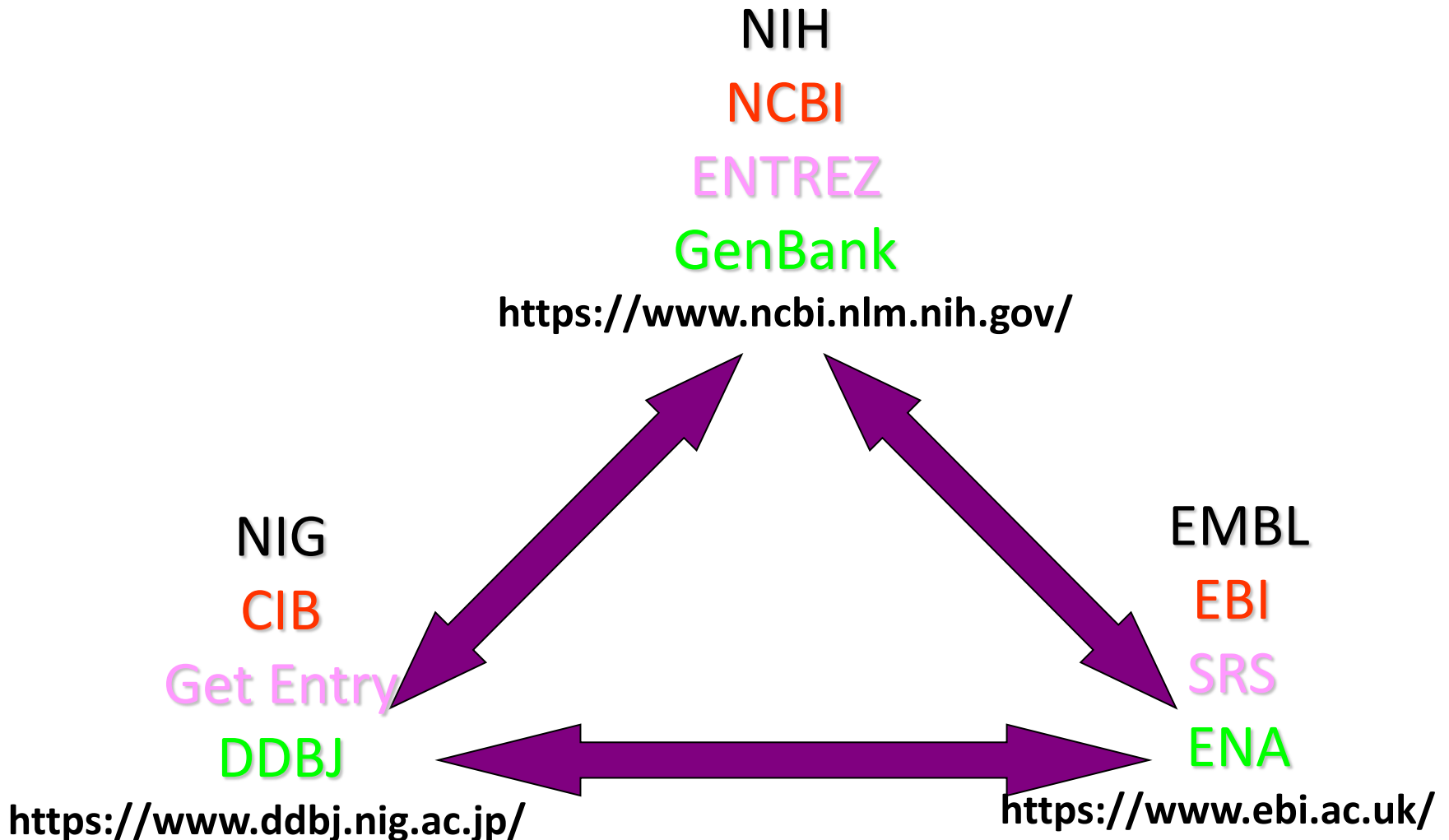
1. Přístup k datům a databáze (lekce - úvod)
2. Zápis sekvencí, hledání podobnosti
3. Konverze dat a formátů
4. Sestavení kompletních sekvencí ze sekvenačních dat (lekce – anotace)
5. Výpočetní analýza sekvencí
6. Návrh oligonukleotidů

Nejdůležitější databáze sekvencí nukleových kyselin a proteinů



- V každém ze tří hlavních bioinformatických center je spravována **genomová databáze** sekvencí nukleových kyselin a odpovídajících, z nich přeložených proteinů.
 - **EMBL Nucleotide Sequence Database / European Nucleotide Archive** (v rámci institutu EBI) – 1980
 - **GenBank** (v rámci institutu NCBI) – 1982
 - **DDBJ** (The DNA Data Bank of Japan) - 1984
- Tři samostatné báze vznikly v důsledku potřeby rychlé dostupnosti databáze sekvencí na jednotlivých kontinentech v době, kdy ještě nebyly rozvinuté vysokorychlostní komunikační sítě.

Mezinárodní spolupráce sekvenčních databází (velká trojka)



Sdílení dat v základních databázích

V každém z bioinformatických center jsou dostupné jednoduché nástroje pro manipulaci s daty



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

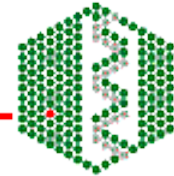
National Center for Biotechnology Information (NCBI)

EMBL: <http://www.ebi.ac.uk>

European Bioinformatics Institute (EBI)

EMBL

European Bioinformatics Institute



DDBJ: <http://www.ddbj.nig.ac.jp/>

National Institute of Genetics (NIG)



ExPASy:

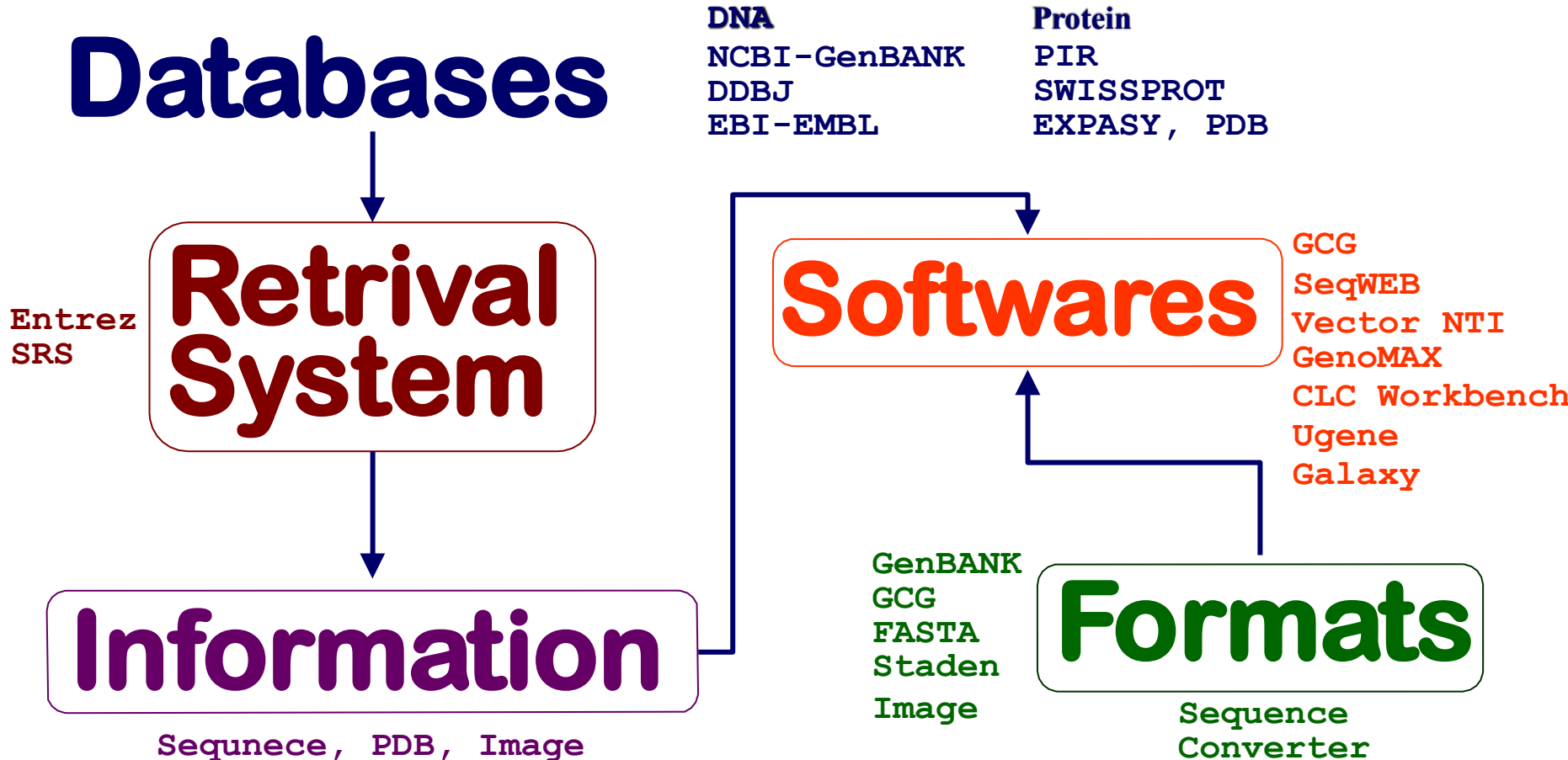
Expert Protein Analysis System



PDB: <https://www.rcsb.org/>

RCSB Protein Data Bank (RCSB PDB)

Získání dat a manipulace se sekvencemi



2. Zápis sekvence



- **Sekvence** – zápis posloupnosti jednoznačných znaků odpovídajících jednotlivým zbytkům (monomerům), které se nacházejí v odpovídající posloupnosti v dané makromolekule
 - ◆ **DNA nebo RNA od 5'-konce k 3'-konci**
 - ◆ 5' CAAACGTCGTCTA 3'
 - ◆ **protein od N-konce k C-konci**
 - ◆ (NH₂-) MKRLSALGPGGLTRR (-COOH)
- používají se jednopísmenové kódy dle pravidel IUPAC

Standardní kódy pro sekvence nukleových kyselin podle IUB/IUPAC



A	adenosin
C	cytidin
G	guanidin
T	thymidin
U	uridin
R	G/A (pu <u>R</u> in)
Y	T/C (p <u>Y</u> rimidin)
K	G/T (nukleosid s <u>K</u> eto skupinou)
M	A/C (nukleosid s a <u>M</u> ino skupinou)
S	G/C (silná = <u>S</u> trong vazba)
W	A/T (slabá = <u>W</u> weak vazba)
<hr/>	
B	G/T/C (not A)
D	G/A/T (not C)
H	A/C/T (not G)
V	G/C/A (not T)
N	A/G/C/T (jakýkoli)
-	mezera (gap) neurčené délky

Využití zápisu s degenerovanými nukleotidy

TACGGT

TATAAT

TATAAT

GATACT

TATGAT

TATATT

Konsenzní sekvence: **TATAAT**

Degenerovaná sekv.: **KAYRNT**

Standardní kódy pro sekvence aminokyselin podle IUB/IUPAC

A	alanin
B	kys. asparagová nebo asparagin
C	cystein
D	kys. asparagová
E	kys. glutamová
F	fenylalanin
G	glycin
H	histidin
I	isoleucin
K	lysin
L	leucin
M	metionin
N	asparagin
P	prolin
Q	glutamin
R	arginin
S	serin
T	treonin
U	selenocystein
V	valin
W	tryptofan
Y	tyrosin
Z	kys. glutamová nebo glutamin
X	jakákoli aminokyselina
*	translační stop (terminační kodon)
-	mezera (gap) neurčené délky

Běžné formáty sekvencí

- Prostý text
- FASTA
- FASTQ
- Genbank
- EMBL
- GCG
- PIR
- ASN1
- Výstupní data sekvenování: ABI, AB1, SCF, SFF, BAM, SAM, FASTF aj.

PLAIN SEQUENCE FORMAT

Obsahuje pouze IUPAC znaky

Obsahuje jedinou sekvenci

Příklad

```
AACCTGCGGAAGGATCATTACCGAGTGCGGGTCCTTTGGGCCCAA  
CCTCCCATCCGTGTCTATTGTAC
```

Použití: pro zápis krátké sekvence např. v textu nebo obrázku

FASTA FORMAT



Může obsahovat více sekvencí

Začíná specifickým záhlavím „>“, za kterým následuje definice

Příklad:

```
>U03518 Aspergillus awamori internal transcribed spacer 1 (ITS1)
AACCTGCGGAAGGATCATTACCGAGTGCGGGTCCTTTGGGCCCAACCTCCCATCCGTGTCTATTGTACCC
TGTTGCTTCGGCGGGCCCGCCGCTTGTCGGCCGCCGGGGGGGCGCCTCTGCCCCCGGGCCCGTGCCCGC
CGGAGACCCCAACACGAACACTGTCTGAAAG
```

```
>LinB_protein
MSLGAKPFGEKKFIEIKGRRMAYIDEGTGDPILFQHGNPTSSYLWRNIMPHCAGLGR
LIACDLIGMGDSDKLDPSGPERYAYA EHRDYLDALWEALDLGDRVVLVVDWGSALG
FDWARRHRERVQGIAYMEA IAMPIEWADFPEQDRDLFQA FR SQAGEELVLQDNVFE
QVLPGLILRPLSEAEMAAYREPFLAAGEARRPTLSWPRQIPIAGTPADVVAIARDYA
GWLSESP I PKLFINAEPGALTTGRMRDFCRTWPNQTEITVAGAHFIQEDSPDEIGAA
IAAFVRRLRPA
```

Použití: univerzální formát pro zápis sekvence vhodný jako vstupní data pro většinu software.



FastQ FORMAT

Záhlaví obsahuje automaticky generovaný identifikátor klastru ze sekvenování.
Následuje primární sekvence
a informaci o kvalitě stanovení sekvence

Příklad:

```
@HWUSI-EAS100R:6:73:941:1973#0/1  
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT  
+  
!''*(((((***+))%%%+))(%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```

Použití: výstupní data při sekvenování nové generace

KLÍČ K IDENTIFIKÁTORŮM KVALITY:

Nejnižší kvalita

nejvyšší kvalita

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\\]^_`abcdefghijklmnopqrstuvwxyz{|}~
```

Viz též Phred Quality Score

GENBANK FORMAT



Začíná řádkem s názvem pole LOCUS

Obsahuje mnoho různých deskriptorů

Začátek primární sekvence je vyznačen ORIGIN a sekvence je ukončena „//“

Zpravidla obsahuje i anotované sekvence proteinů a odkazy do dalších databází

Příklad:

```
LOCUS      AAU03518      237 bp      DNA              PLN              04-FEB-1995
DEFINITION Aspergillus awamori internal transcribed spacer 1 (ITS1) and 18S
           rRNA and 5.8S rRNA genes, partial sequence.
ACCESSION  U03518
VERSION    U03518.1  GI 1235658
BASE COUNT      41 a      77 c      67 g      52 t
ORIGIN
      1  aacctgcgga  aggatcatta  ccgagtgcgg  gtcctttggg  cccaacctcc  catccgtgtc
     61  tattgtacc  tgttgcttcg  gcgggcccg  cgcttgtcgg  ccgccggggg  ggcgccctctg
    121  cccccgggc  ccgtgcccg  cgagacccc  aacacgaaca  ctgtctgaaa  gcgtgcagtc
    181  tgagttgatt  gaatgcaatc  agttaaact  ttcaacaatg  gatctcttgg  ttccggc
//
```

EMBL FORMAT



Začíná řádkem s jedinečným identifikátorem (ID), následuje anotace.

Obsahuje mnoho různých deskriptorů

Sekvence začíná symboly SQ a sekvence je ukončena „//“

Zpravidla obsahuje i anotované sekvence proteinů a odkazy do dalších databází

Příklad:

```
ID   AA03518      standard; DNA; FUN; 237 BP.
XX
AC   U03518;
XX
DE   Aspergillus awamori internal transcribed spacer 1 (ITS1) and 18S
DE   rRNA and 5.8S rRNA genes, partial sequence.
XX
SQ   Sequence 237 BP; 41 A; 77 C; 67 G; 52 T; 0 other;
      aacctgcgga aggatcatta ccgagtgcgg gtcctttggg cccaacctcc catccgtgtc      60
      tattgtaccc tgttgcttcg gcgggcccgc cgcttgtcgg ccgccggggg ggcgcctctg      120
      cccccgggc ccgtgccgc cggagacccc aacacgaaca ctgtctgaaa gcgtgcagtc      180
      tgagttgatt gaatgcaatc agttaaact ttcaacaatg gatctcttgg ttccggc      237
//
```

Formáty sekvencí obsahující mnohonásobná příložená

- Sekvenční příložená umožňuje srovnat podobné sekvence
- K dispozici celá řada formátů, obvykle přizpůsobených používaným programům
 - Multi FASTA
 - Phylip
 - PAUP / NEXUS
 - Clustal
 - MSF



CLUSTAL/MUSCLE MULTIPLE FORMAT

Začíná řádkem s definicí

Vkládá mezery do sekvence tak, aby při mnohonásobném přiložení byly identické zbytky nad sebou

Konzervované pozice se stejným typem zbytku označuje na posledním řádku hvězdičkou

Podobné typy zbytků mohou některé formáty znázorňovat dvoutečkou nebo tečkou

Příklad:

```
Moorella_thermoacetica_ATCC_39073_-_rna.40      -----GTTTGATCCTGGCTCAGGACAAACGCTGGCGGCGTGCCTAACACATGCAA 50
Ammonifex_degensii_KC4_-_rna.5                 -----AGGGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCCTAACACATGCAA 53
Ammonifex_degensii_KC4_-_rna.31                -----AGGGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCCTAACACATGCAA 53
Candidatus_Desulforudis_audaxviator_MP104C_-_DAUD_RS00700 TTTATGGAGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCCTAACACATGCAA 60
Candidatus_Desulforudis_audaxviator_MP104C_-_DAUD_RS06920 TTTATGGAGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCCTAACACATGCAA 60
*****:*****:*****

Moorella_thermoacetica_ATCC_39073_-_rna.40      GTCGAGCGGTCCTTAATTGGGAAATCTTCGGATGGAACCGATTAAGATAGCGGCGGAC 110
Ammonifex_degensii_KC4_-_rna.5                 GTCGAGCGGGCTT-----GTCAGGGCCTTGTGT----CCTGGCAAGTTGAGCGGCGGAC 103
Ammonifex_degensii_KC4_-_rna.31                GTCGAGCGGGCTT-----GTCAGGGCCTTGTGT----CCTGGCAAGTTGAGCGGCGGAC 103
Candidatus_Desulforudis_audaxviator_MP104C_-_DAUD_RS00700 GTCGTGCGA---TTGAGAGGTGAGCATCTCACTT----CTCAA-----GAGCGGCGGAC 107
Candidatus_Desulforudis_audaxviator_MP104C_-_DAUD_RS06920 GTCGTGCGA---TTGAGAGGTGAGCATCTCACTT----CTCAA-----GAGCGGCGGAC 107
**** *: ***** ** *: :*: * *: :*: *****
```

Poznámka k používaným fontům

■ Proporcionální fonty

- ◆ Arial, Times
- ◆ Každý znak - jiná šířka
- ◆ Nevhodné pro zápis sekvence

gaattttttt
cttaaaaaaa

■ Neproporcionální fonty

- ◆ Vhodné pro zápis sekvence
- ◆ Všechny znaky stejná šířka
- ◆ Courier, Monospaced

gaatttttttt
cttaaaaaaaa

- K editaci jsou vhodné editory, které neukládají informace o formátu textu (Notepad, vývojářské editory – PSPad, aj.)

- Některé formáty jako např. GCG obsahují vnitřní kontrolní součty

Surová data – elektroforetogramy ze sekvenování v kapiláře

Různé formáty

- ◆ *.abi
- ◆ *.ab1
- ◆ *.scf

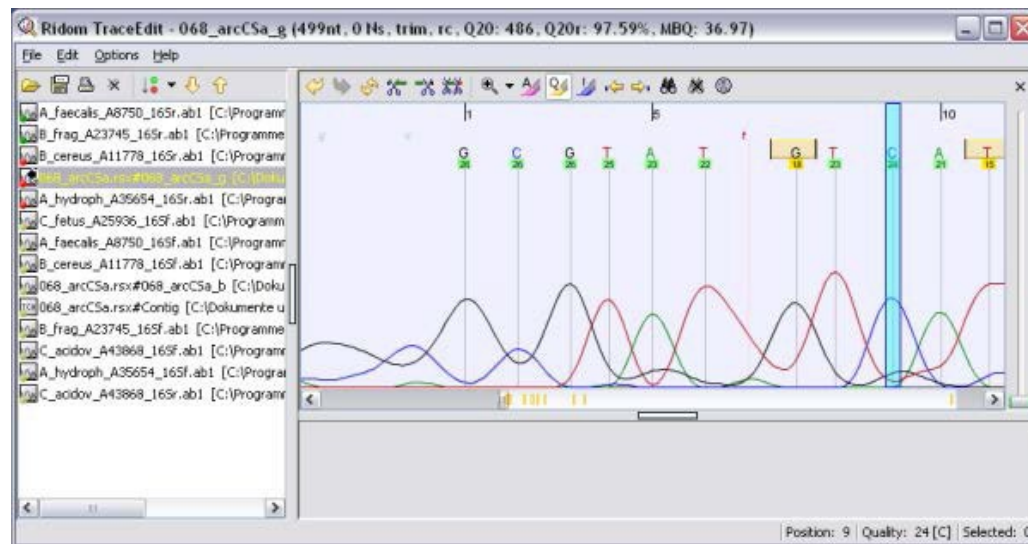
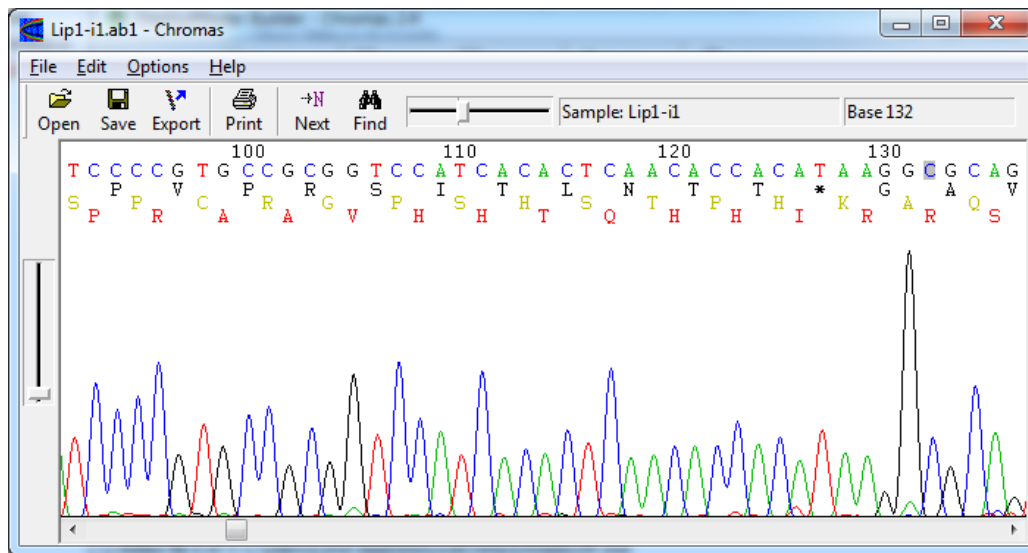
Prohlížeče, např.

- ◆ Chromas Lite
- ◆ ABIView
- ◆ Ridom Trace Edit

Export

- ◆ FASTA
- ◆ Prostý text

- Formáty z NGS vyžadují složitější zpracování



Jednoduché formáty sekvencí mají omezení a neobsahují

- Data o expresi genů
- Variace a polymorfismy
- Specifické informace o zdroji sekvence (organismu, klonech, ...)
- Odkazy na další informace
- Informace o kvalitě

3. Konverze dat a hledání motivů

- **Převod informace mezi řetězci**
 - Reverse-complement
- **Hledání motivů**
 - Přesné
 - Podobné
- **Přepis a překlad podle ústředního dogmatu**
 - Transkripce
 - Translace – genetický kód
- **Sekvenční příložen**
 - Párové, stanovení identity a podobnosti
 - Mnohonásobné, identifikace konzervativních motivů
- **Spojování, rozdělování**
 - Restrikční štěpění
 - Klonování *in silico*, konstrukce vektorů a rekombinantní DNA pro přípravu proteinů
- **Assembly – kompletace a sestavení genomů**

Příklady nástrojů pro konverzi formátů

■ UNIX-GCG

- To Genbank, To Fasta....
- From Genbank, From Fasta...

■ SEQRET

- <https://www.ebi.ac.uk/Tools/sfc/>

■ SMS – The Sequence Manipulation Suite v2

- ◆ <http://www.bioinformatics.org/sms2/>

- EMBL to FASTA
- GenBank to FASTA
- Reverse Complement – převod mezi řetězci
- Filtrování znaků

■ Vzájemná konverze - Sanger, IonTorrent, Illumina, (ONT)

- Biopython
- EMBOSS
- BioPerl
- Samtools

Převod informace mezi řetězci



Nástroj Reverse Complement

http://www.bioinformatics.org/sms2/rev_comp.html

Převod mezi dvěma řetězci (pozitivní-kódující-horní/negativní-antikódující-spodní)

5' CCCCATGTTT 3'

3' GGGGTACAAA 5'

>Sample sequence

5' CCCCATGTTT 3'

>Reverse complement **dle pravidel o párování bází**

5' AAACATGGGG 3'

>Reverse - nemá biologický význam

5' TTTGTACCCC 3'

>Complement - nemá biologický význam

5' GGGGTACAAA 3'

Hledání motivů v sekvencích

Hledání slov = uspořádaná množina znaků


GAATTC

GARYTC

GAAN (1-50) **TTC**

- Přesně definovaný motiv
- Degenerované symboly
- Povoleno počet neshod
- Motivy od sebe vzdálené x zbytků

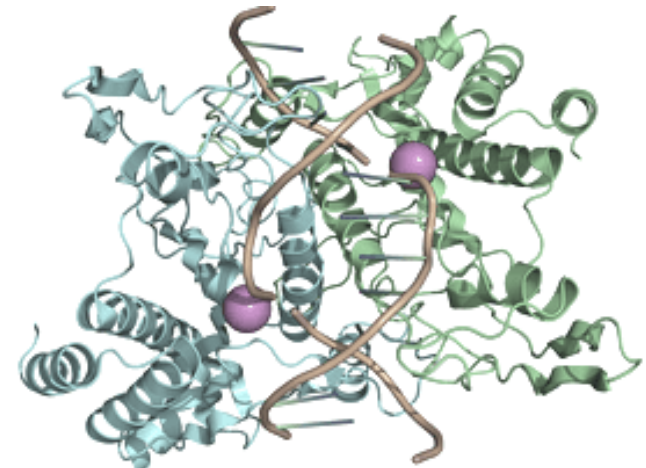
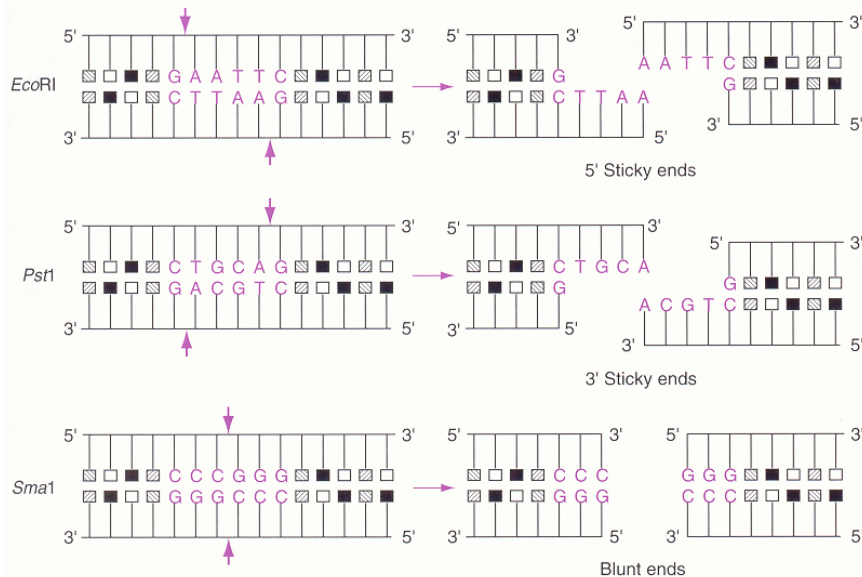
Standardní příklady hledání

- Restrikční místa
 - Repetice
 - ◆ Přímé
 - ◆ Obrácené (vlásenky se smyčkou)
 - Konsenzní vzory
 - Uživatelem definované vzory
 - Otevřené čtecí rámce (START STOP)
- 
- Základ pro hledání genů a funkčních oblastí

Restrikční endonukleázy třídy II mají praktické využití



- Vážou se na specifické (4-6 pb) sekvence nukleotidů
- Katalyzují štěpení dvou řetězců molekuly DNA uvnitř vazebného místa nebo v jeho bezprostředním sousedství
- Produkty štěpení RE
 - tupé konce (po štěpení obou řetězců ve stejném místě)
 - přečnívající konce (po štěpení řetězců v různých místech, která jsou obvykle vzdálena 1-4 nukleotidy)
 - - 5' přečnívající
 - - 3' přečnívající



Restrikční analýza *in silico*

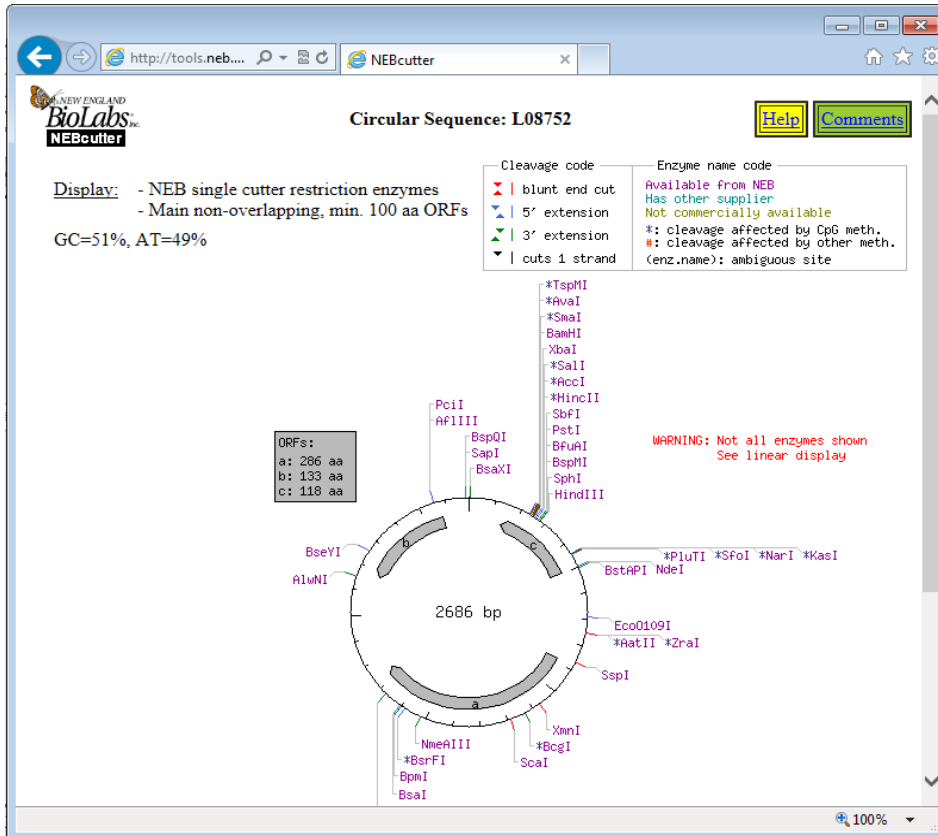
- Restrikční endonukleázy třídy II
 - ◆ Sekvenčně specifické endonukleázy, které štěpí DNA v rozpoznávaných sekvencích
 - ◆ Přehled dostupný v databázi REBASE- Restriction Enzyme Database
<http://rebase.neb.com/rebase/rebase.html>
 - ◆ Sekvence rozpoznávacích míst
 - ◆ Producent enzymu
 - ◆ Reference
 - ◆ Komerční dostupnost
 - ◆ Sekvence genů
 - ◆ Krystalografická data
 - ◆ Citlivost k metylaci
 - ◆ REBpredictor – predikce rozpoznávací sekvence u nových enzymů
 - ◆ Rebase genomes – identifikace genů pro RE v genomech

Software pro restriční mapování

- Provádí hledání restričních míst na základě analýzy sekvence DNA
- Konstrukce restričních map
 - ◆ Nezbytný předpoklad pro klonování
 - ◆ Interpretace RFLP polymorfizmů
 - ◆ Simulace výsledků gelové elektroforézy restričních fragmentů
- Virtuální klonování
- Vytvoření kvalitní grafiky ilustrující restriční mapy
 - ◆ RestrictionMapper (<http://www.restrictionmapper.org/>)
 - ◆ WebCutter
 - ◆ NEB Cutter v3.0 (<https://nc3.neb.com/NEBcutter/>)
 - ◆ EMBOSS Restrict (<https://www.bioinformatics.nl/cgi-bin/emboss/restrict>)
 - ◆ pDRAW32 freeware (<http://www.acaclone.com/>)

NEB Cutter

<http://tools.neb.com/NEBcutter2/>



- Enzymy – výstup tabulka
 - ◆ kompletní sada
 - ◆ komerční sada
 - ◆ které sekvenci neštěpí
 - ◆ které štěpí – počet a pozice rozpoznávacích míst
- Lineární nebo kružnicová mapa sekvence se znázorněním pozice restričních míst
 - ◆ Grafika
 - ◆ Identifikace ORF a translace do proteinu

Vyhledání otevřených čtecích rámců



- **ORF (Open Reading Frame)**
Sada překládaných kodonů mezi iniciačním a terminačním kodonem
- Výsledek je závislý na použitém genetickém kódu
 - Databáze genetických kódů v NCBI
 - <https://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi>
- U prokaryot, které nemají introny je základem hledání genů
- U eukaryot zpravidla využíváme analýzu ORF u komplementární DNA (cDNA) vzniklé reverzní transkripcí z mRNA

ORF Finder (Open Reading Frame Finder)

<https://www.ncbi.nlm.nih.gov/orffinder/>

Open Reading Frame Finder

ORF finder searches for open reading frames (ORFs) in the DNA sequence you enter. The program returns the range of each ORF, along with its protein translation. Use ORF finder to search newly sequenced DNA for potential protein encoding segments, verify predicted protein using newly developed SMART BLAST or regular BLASTP.

This web version of the ORF finder is limited to the subrange of the query sequence up to 50 kb long. Stand-alone version, which doesn't have query sequence length limitation, is available for [Linux x64](#).

Examples (click to set values, then click Submit button) :

- NC_011604 Salmonella enterica plasmid pWES-1; genetic code: 11; 'ATG' and alternative initiation codons; minimal ORF length: 300 nt
- NM_000059; genetic code: 1; start codon: 'ATG only'; minimal ORF length: 150 nt



Enter Query Sequence

Enter accession number, gi, or nucleotide sequence in FASTA format:

```
>bacteriophage 3A
TCGCTTAAACCTTCATGCCTTCTGGACACCTAAATGGTCTAATTCAGCTCCAABGGTCATGCCTTCTACTT
TTCATATTAACCTCCCTTCTAGCTTCCAAAAAGTT@TCTTAAATCCGTACCTGTAATGACTTTTGTTCACCTT
TCTTCAGTCTCTT@CTTTATTCTCTTCAATTAAGTATTTCTAAAAGTTTACATACGGCTGTTTTCTGACTTCAG
GTCCACCATACTGCTCCATACAGAAACGTTGATTTTCTTAATGTTTCG@ATAAAATATCTTTATTGAGATTGT
TGCTTTCCCATCTCTCTGGTTCAGT@TCTGAATCTTCTCATCTTCACCATTGATTTCTCGAAATATATCTT
GCTTT@TGATAAGTTTTAGTGCTCATCTTGTCAAAACATCTTCTCAGTCAATCCTTCATC@TTTAAATAAT
AATAACTGTCGCTTTTTGTCTCATTTTTGTTGCGTTAGG@TACTTCTTTTTATTCTCTTGATTACTAATTC
```

From: To:

Choose Search Parameters

- Minimal ORF length (nt):
- Genetic code:
- ORF start codon to use:
 - "ATG" only
 - "ATG" and alternative initiation codons
 - Any sense codon
- Ignore nested ORFs:

Start Search / Clear

ORF Finder (Open Reading Frame Finder)

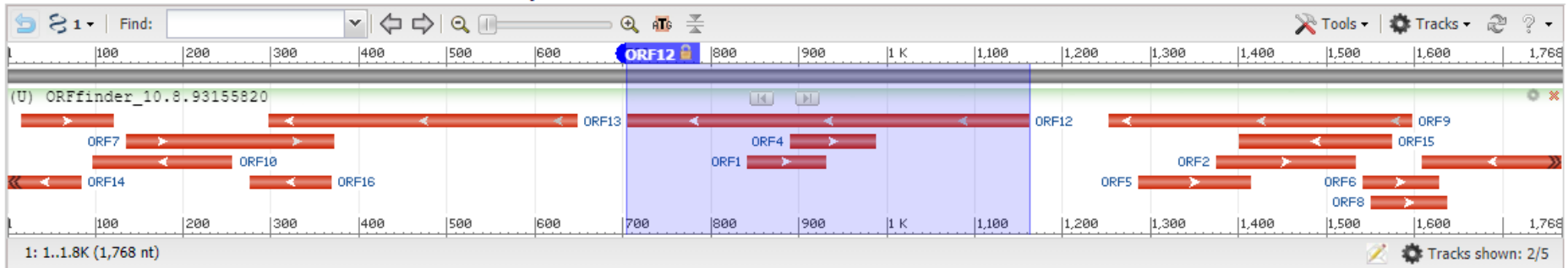
<https://www.ncbi.nlm.nih.gov/orffinder/>

Open Reading Frame Viewer

Help

Sequence

ORFs found: 16 Genetic code: 11 Start codon: 'ATG' only



Six-frame translation...

ORF12 (151 aa)

Display ORF as...

Mark

Mark subset...

Marked: 0

Download marked set

as Protein FASTA

```
>|c1|ORF12
MTKTLKVKYKGGDQVASEQEGGKVSVTLSNLEADTTYPKGTQVVAEENGK
ESSKVDVPQFKTNPILVSGVSFTPETKSITVNADDNVEPNIAPISTATNKT
LKYTSEHPFVTVDERTGAIHGVAEGTSVITATSTDGSOKSGQITVTVTN
G
```

Label	Strand	Frame	Start	Stop	Length (nt aa)
ORF12	-	2	1161	706	456 151
ORF13	-	2	648	298	351 116
ORF9	-	1	1597	1253	345 114
ORF7	+	3	135	371	237 78
ORF15	-	3	1574	1401	174 57
ORF2	+	1	1375	1533	159 52
ORF10	-	1	256	98	159 52
ORF11	-	2	<1767	1609	159 52
ORF5	+	2	1286	1414	129 42
ORF3	+	2	17	121	105 34

ORF12

Marked set (0)

SmartBLAST

SmartBLAST best hit titles...

BLAST

BLAST

BLAST Database:

UniProtKB/Swiss-Prot (swissprot)

Translace *in silico*



Překlad genetické informace z DNA do proteinu

- 6 možných čtecích rámců ve dvouřetězcové DNA
- Vymezené oblasti - exony
- Výběr genetického kódu
- Nástroje: EMBOSS Transeq, EMBOSS Sixpack, aj.

```
E F K T S K S C E K A I T K * R * F G Y F1
N S K P A K A V K K P L P S K D N L A I F2
I Q N Q Q K L * K S H Y Q V K I I W L Y F3
1 GAATTCAAAAACAGCAAAAAGCTGTGAAAAAGCCATTACCAAGTAAAGATAATTTGGCTAT 60
----:----|----:----|----:----|----:----|----:----|----:----|
1 CTTAAGTTTTGGTCGTTTTTCGACACTTTTTCGGTAATGGTTCATTTCTATTAACCGATA 60
S N L V L L L Q S F A M V L Y L Y N P * F6
X I * F W C F S H F L W * W T F I I Q S F5
F E F G A F A T F F G N G L L S L K A I F4
```

EMBOSS Transeq

http://www.ebi.ac.uk/Tools/st/emboss_transeq/

EMBL-EBI [Services](#) [Research](#) [Training](#) [Industry](#) [About us](#)

EMBOSS Transeq

[Input form](#) [Web services](#) [Help & Documentation](#) [Share](#) [Feedback](#)

[Tools](#) > [Sequence Translation](#) > EMBOSS Transeq

EMBOSS Transeq

EMBOSS Transeq translates nucleic acid sequences to their corresponding peptide sequences. It can translate to the three forward and three reverse frames, and output multiple frame translations at once.

STEP 1 - Enter your input sequence

Enter or paste a set of sequences in any supported format:

Or, [upload](#) a file:

1
2
3
F (Forward three frames)
-1
-2
-3
R (Reverse three frames)
6 (All six frames)

CODON TABLE

ost users and, for that reason, are not visible.

100%

Příklady translace *in silico*



ExPASy
Bioinformatics Resource Portal

Translate Tool - Results of translation

Open reading frames are highlighted in **red**. Please select one of the following frames

5'3' Frame 1

```
LLIQQAKSNSDTPAMPLDTCGAMSQGMIGYWLETEINRILTEMNSDRTVGTIVTRVEVD  
KDDPRFDNPTKPIGPFYTKEEVEELQKEQPDSVFKEDAGRGYRQVVASPLPQSILEHQLI  
QTLADGKNIVIACGGGGIPVIKKENTYEGVEA
```

5'3' Frame 2

```
Y-SNKLNRVTQRRQCHWILVVQCHRV--AIGWKLKSI AF-LK-IVIEL-AQSLHVWK-I  
KMIHDLITQLNQLVLFIRKKKLKNYKKNSQTQSLKMKQDVVIEK-LRHHYLNLY-NTS-F  
KL-QTVKILSLHAVVAVFQL-KKKIPMKVLK
```

5'3' Frame 3

```
INPTS-IEQ-HNAGNAIGYLWCNVTGYDRLLVGN-NQSHFN-NE---NCRHNR YTCGSR-  
R-STI--PN-TNWSFLYERRS-RITKRTARLSL-RRCRTWL-KSSCVTTTTSIYTRTPVNS  
NFSRR-KYCHC MRWWRYSYKKRKYL-RC-S
```


Klonování *in silico*, konstrukce vektorů

- Kombinace segmentů sekvencí
 - ◆ známé/neznámé funkce
- Plazmidy
 - ◆ přebírané z databáze
 - ◆ zpravidla známé funkce
 - ◆ regulační sekvence pro expresi
- Inzerty – obvykle nové sekvence
 - ◆ charakterizované restriční mapou
 - ◆ charakterizované sekvencí DNA
 - ◆ charakterizované funkcí
- Design *in vitro* mutageneze
- Nomenklatura pro konstrukty není stanovena

Clone Manager (Sci-Ed Software)

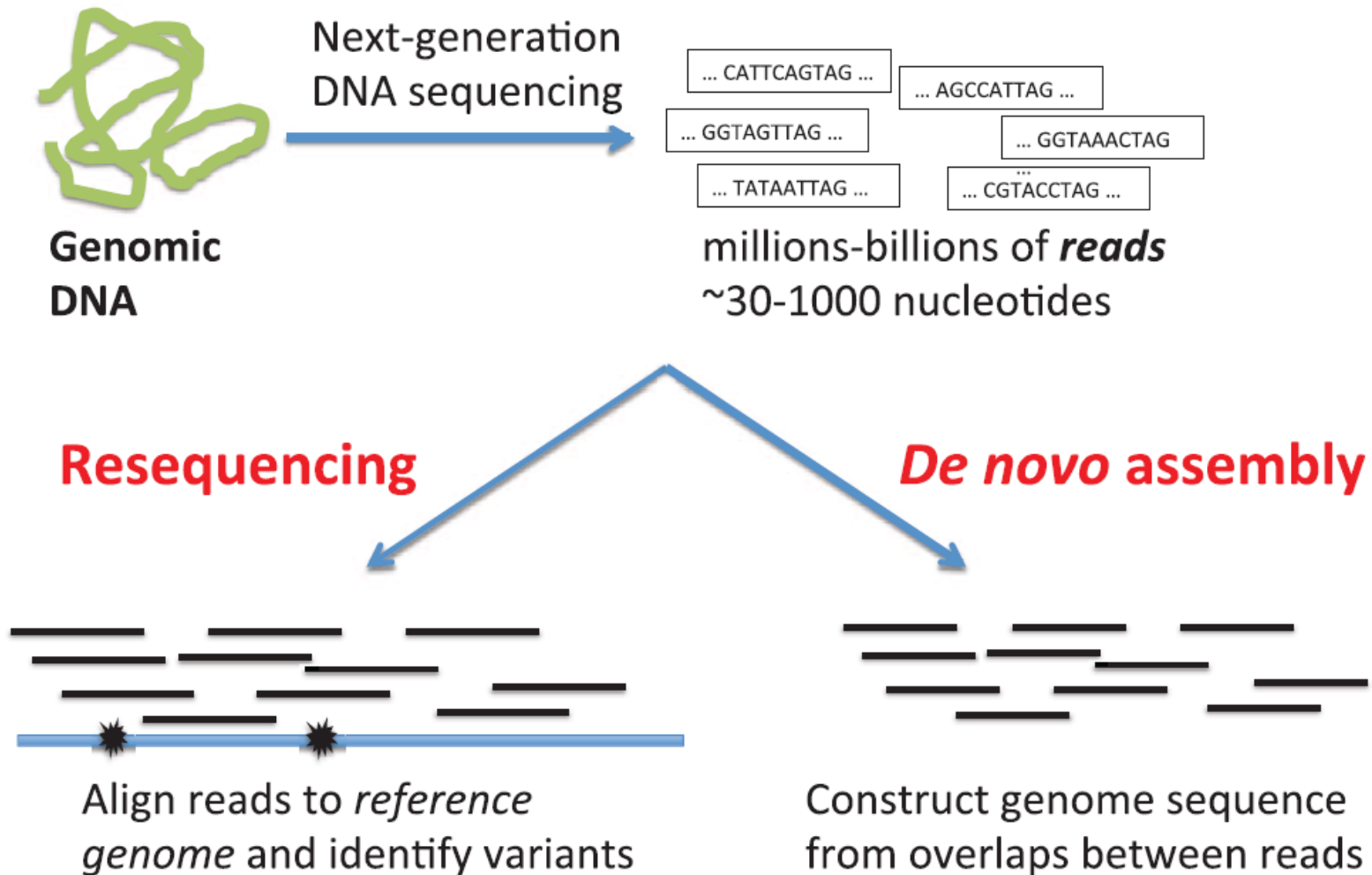
http://www.sci-ed.com/pr_cmbas.htm

The screenshot displays the Clone Manager software interface. The main window shows a circular plasmid map for a 2686bp construct named SYNPU18V. The map is annotated with various restriction enzyme sites. A list of 44 enzyme sites is displayed on the right side of the interface, with the ApoI site highlighted in blue. The list includes the enzyme name, its position (Pos), and its type (Type).

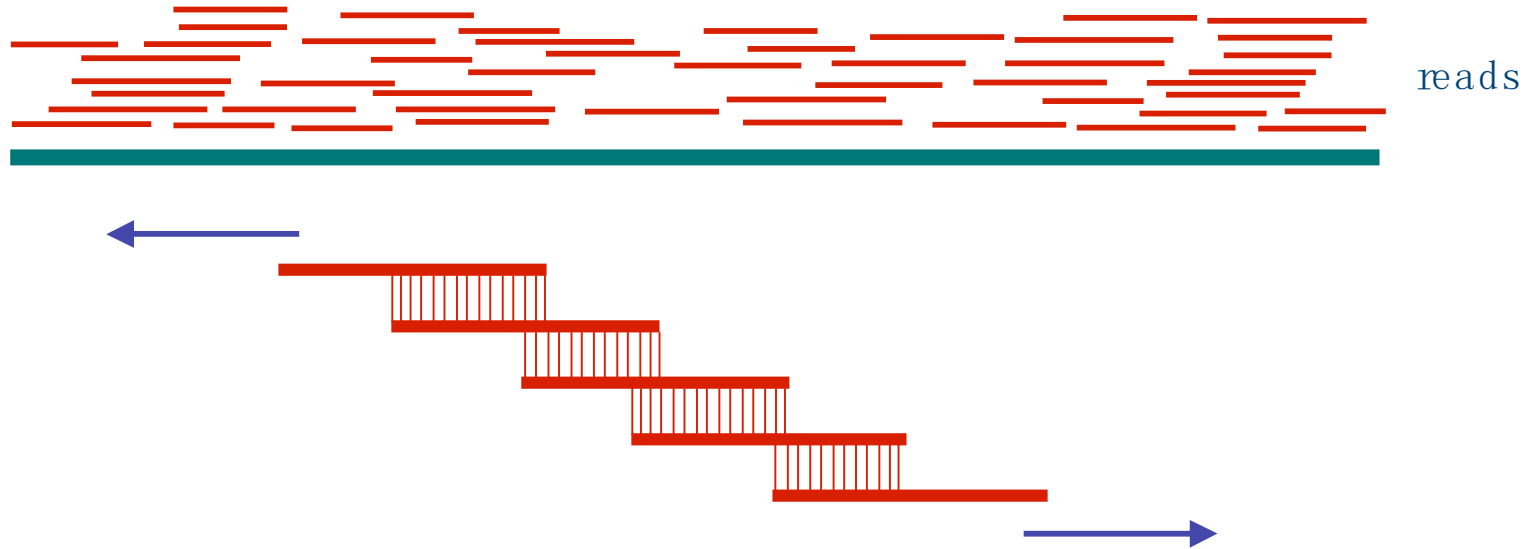
Name	Pos	Type
ApoI	230	sc 5'
EcoRI	230	sc 5'
BanII	236	sc 3'
Eco53kI	236	sc bl
SacI	236	sc 3'
Acc65I	242	sc 5'
KpnI	242	sc 3'
AvaI	246	sc 5'
SmaI	246	sc bl
XmaI	246	sc 5'
BamHI	251	sc 5'
XbaI	257	sc 5'
AccI	263	sc 5'
HincII	263	sc bl
SalI	263	sc 5'
BspMI	267	sc 5'
SbfI	268	sc 3'
PstI	269	sc 3'
SphI	275	sc 3'

Assembly/ kompletace a sestavení

Resekvenování vs. *de novo* sekvenování



Princip assembly



Pokrytí oblastí $>x$ -násobnou redundancí

Identifikace překryvů, sekvenční příložen

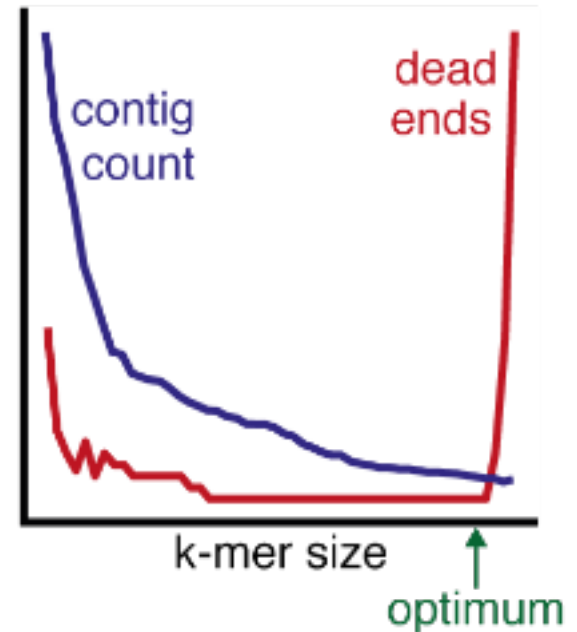
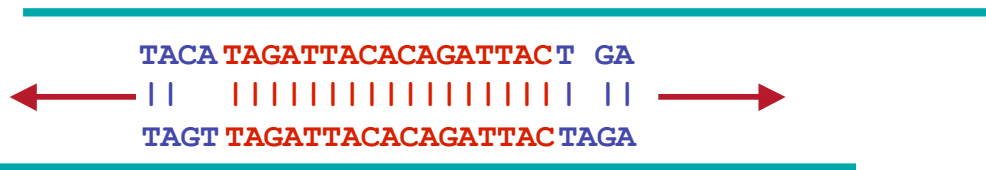
a rekonstrukce sekvence

De novo assembly

- Nezávislé na referenčním genomu
- Parametry
 - Délka čtení
 - Pokrytí genomu (coverage)
- Velké množství dostupných algoritmů
 - Znakové metody
 - Grafové metody
- Výpočetně náročné
- Zpravidla vyžaduje optimalizaci pro každou platformu sekvenování

Princip hledání překryvů

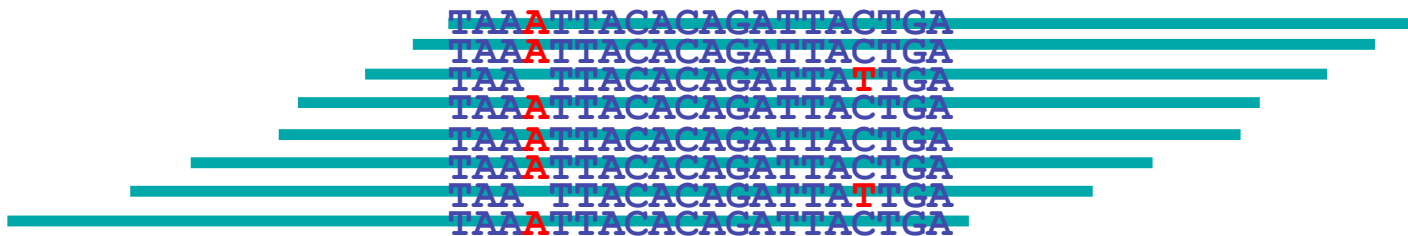
- Vytvoření všech k -merů ve čteních, (např. $k \sim 24$)
- Roztřídění čtení do skupin, které sdílejí k -mer
- Přiložení párů, které sdílejí k -mer
- Mapování a rozšíření sekvenčních příložen



Mapování jednotlivých čtení k referenci



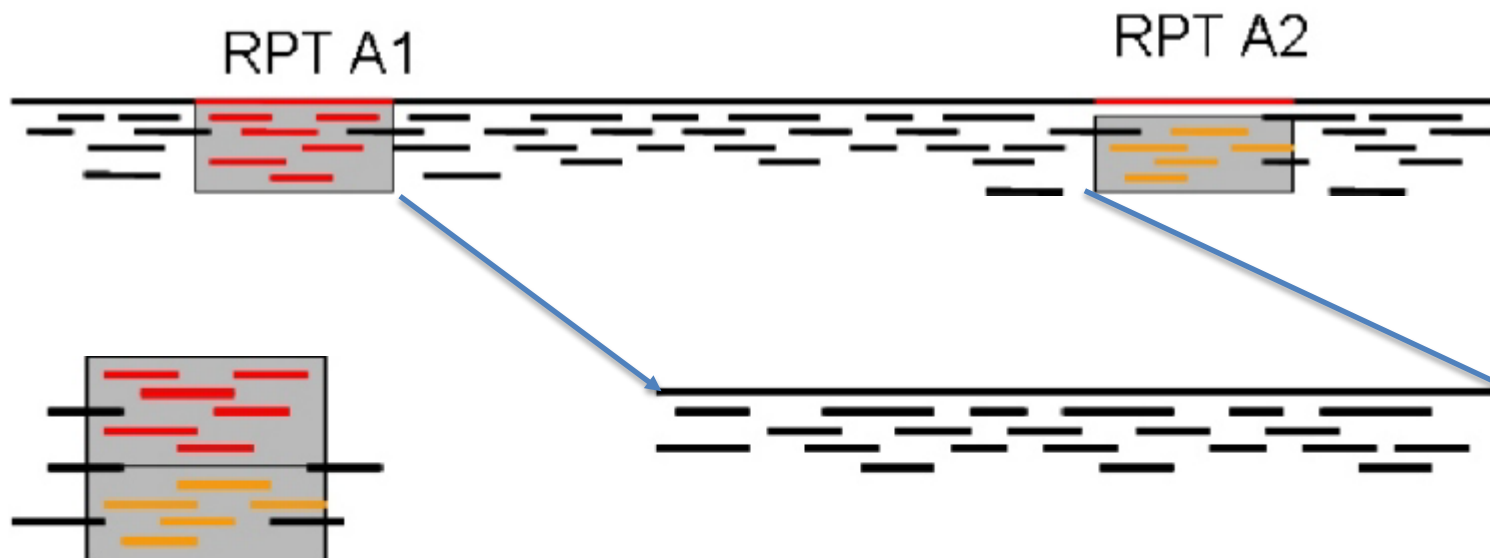
Vytvoření sekvenčního příložení z jednotlivých čtení



- Hloubka pročtení sekvence odráží kvalitu
- Umožňuje vyřešit neshody
 - Chybně stanovené báze
 - Homopolymerní oblasti
- Umožňuje kvantifikovat polymorfizmy

Repetice jsou příčinou rozdělení genomů do kontigů

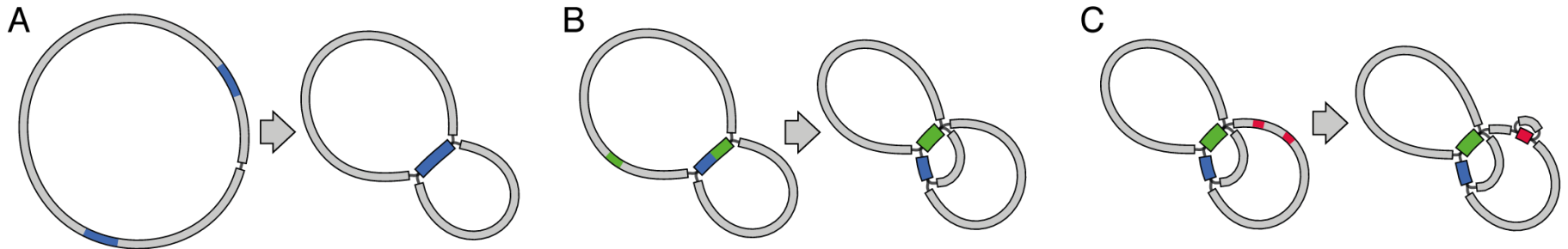
Jestliže čtení je kratší než repetice → nemožnost sestavení sekvence



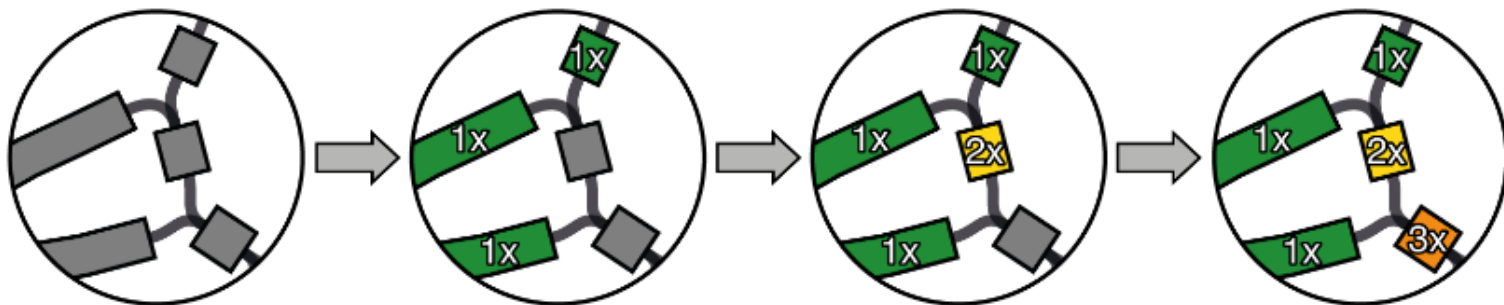
Čtení z mnoha podobných repetitív vedou k vytvoření kontigů s pozměněnou strukturou

Kontig tvořený jedinečnou sekvencí, ohraničený repetitivními sekvencemi

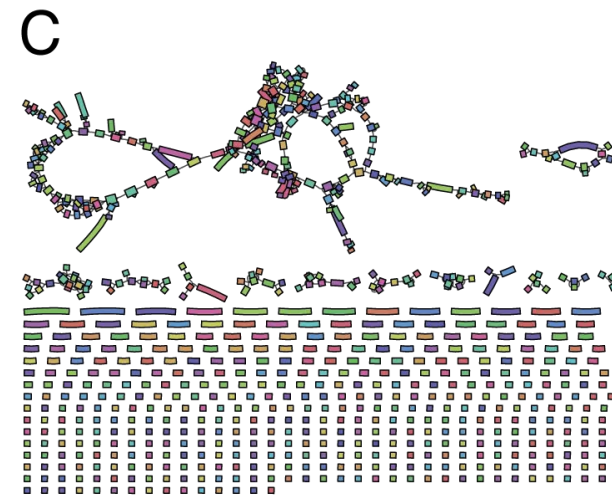
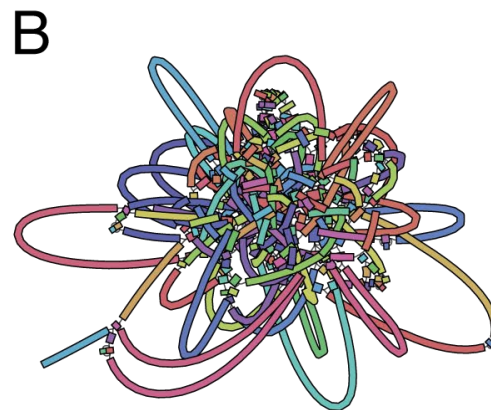
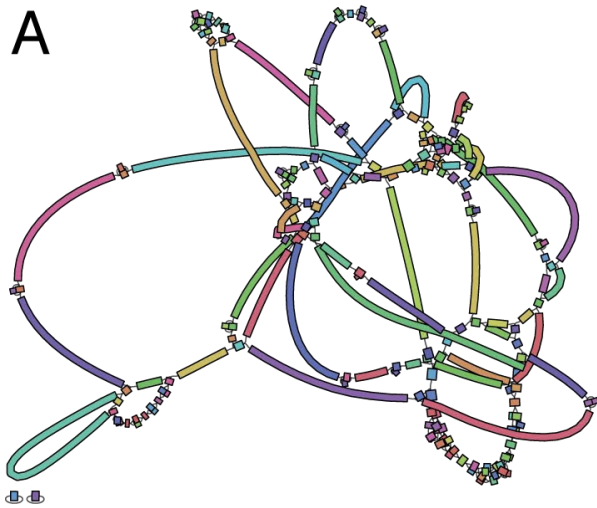
Znázornění repetic v de Bruijnově grafu



- Krátká čtení, hlavní příčina omezení kompletního sestavení
- Stejná sekvence se vyskytuje v genomu vícekrát
- Délka čtení není schopna překlenout tuto repetici
- Pokrytí může indikovat multiplicitu



Příklad de Bruijnova grafu u mikrobiálního genomu (Illumina)



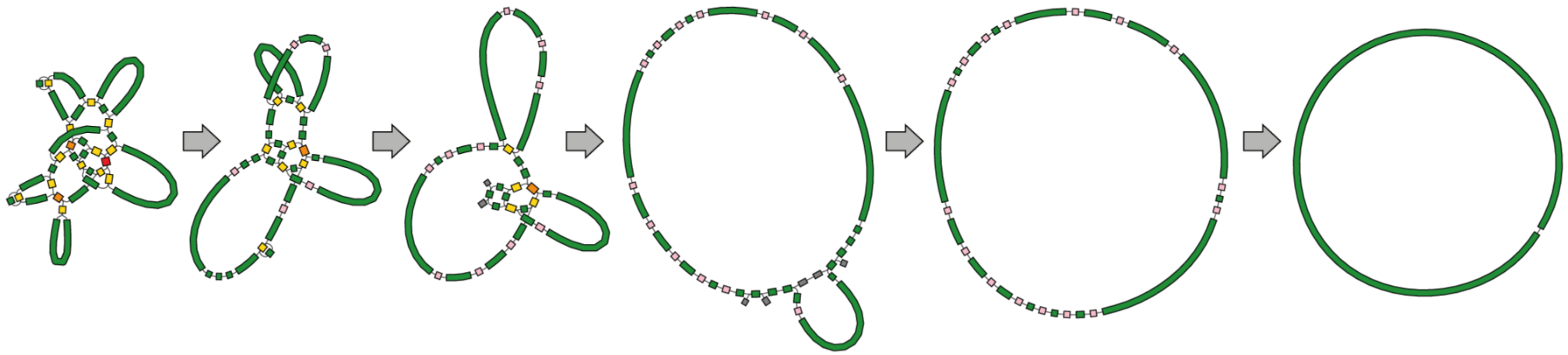
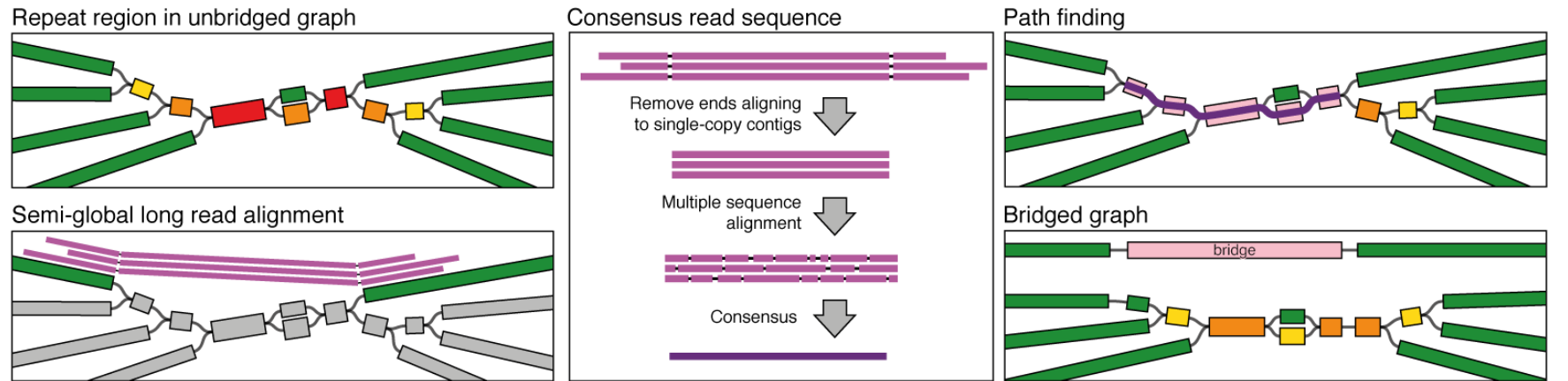
A – kvalitní sestavení

B – sestavení vyžadující optimalizaci, kombinace dlouhých a krátkých kontigů

C – nekvalitní sestavení vycházející z nekvalitních dat, velké množství nezařazených krátkých kontigů

Hybridní assembly a bridging

- Kombinace krátkých čtení (Illumina, IonTorrent) a dlouhých čtení (PacBio, Nanopore) umožňuje hybridní assembly
- Dlouhá čtení: hledání cesty mezi repeticemi



4. Výpočetní analýza sekvencí

- Počet residuí
- Frekvence residuí
- Frekvence oligonukleotidů
- Analýza využití kodonů
- Design oligonukleotidů a primerů

Analýza využití kodonů (codon usage)



- Využití synonymních kodonů
 - ◆ není náhodné
 - ◆ je rozdílné u různých genomů, které mají určité preferované kodony pro určité aminokyseliny
 - ◆ může být problémem při expresi rekombinantních proteinů

- Databáze využití kodonů
<http://www.kazusa.or.jp/codon/>

The Human Codon Usage Table

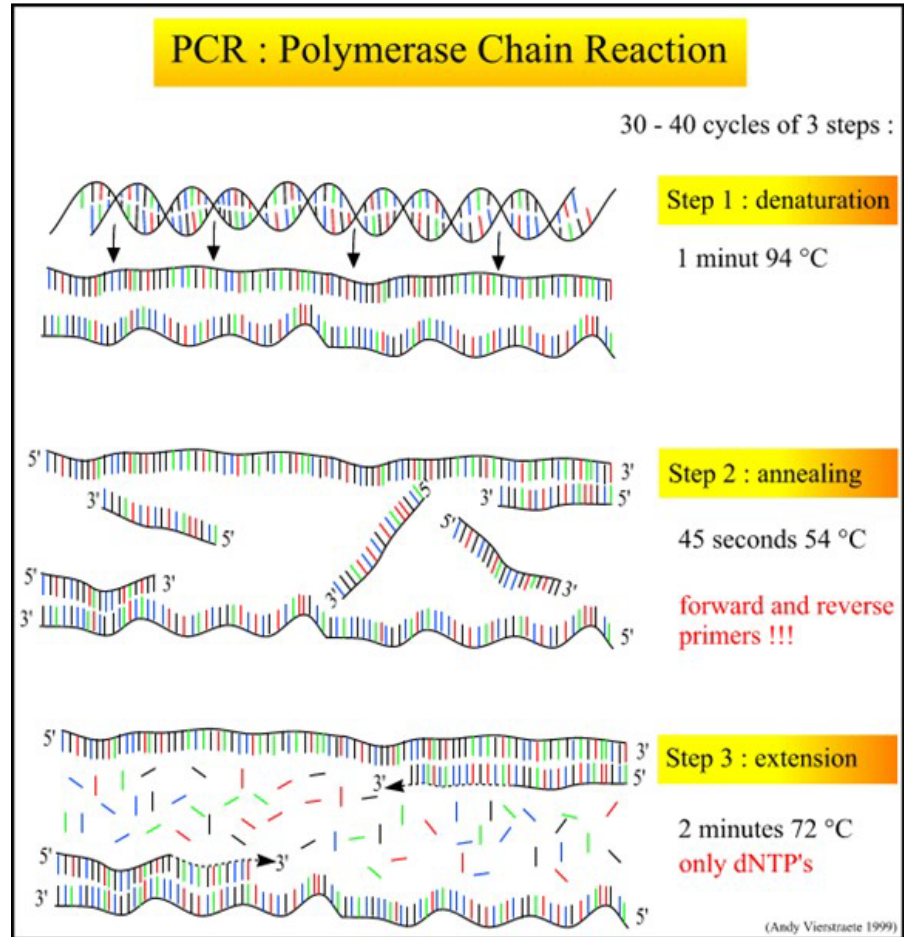
Gly	GGG	17.08	0.23	Arg	AGG	12.09	0.22	Trp	TGG	14.74	1.00	Arg	CGG	10.40	0.19
Gly	GGA	19.31	0.26	Arg	AGA	11.73	0.21	End	TGA	2.64	0.61	Arg	CGA	5.63	0.10
Gly	GGT	13.66	0.18	Ser	AGT	10.18	0.14	Cys	TGT	9.99	0.42	Arg	CGT	5.16	0.09
Gly	GGC	24.94	0.33	Ser	AGC	18.54	0.23	Cys	TGC	13.86	0.58	Arg	CGC	10.82	0.19
Glu	GAG	38.82	0.59	Lys	AAG	33.79	0.60	End	TAG	0.73	0.17	Gln	CAG	32.95	0.73
Glu	GAA	27.51	0.41	Lys	AAA	22.32	0.40	End	TAA	0.95	0.22	Gln	CAA	11.94	0.27
Asp	GAT	21.45	0.44	Asn	AAT	16.43	0.44	Tyr	TAT	11.80	0.42	His	CAT	9.56	0.41
Asp	GAC	27.06	0.56	Asn	AAC	21.30	0.56	Tyr	TAC	16.48	0.58	His	CAC	14.00	0.59
Val	GTC	28.60	0.48	Met	ATG	21.86	1.00	Leu	TTC	11.43	0.12	Leu	CTC	39.93	0.43
Val	GTA	6.09	0.10	Ile	ATA	6.05	0.14	Leu	TTA	5.55	0.06	Leu	CTA	6.42	0.07
Val	GTT	10.30	0.17	Ile	ATT	15.03	0.35	Phe	TTT	15.36	0.43	Leu	CTT	11.24	0.12
Val	GTC	15.01	0.25	Ile	ATC	22.47	0.52	Phe	TTC	20.72	0.57	Leu	CTC	19.14	0.20
Ala	GCG	7.27	0.10	Thr	ACG	6.80	0.12	Ser	TCG	4.38	0.06	Pro	CCG	7.02	0.11
Ala	GCA	15.50	0.22	Thr	ACA	15.04	0.27	Ser	TCA	10.96	0.15	Pro	CCA	17.11	0.27
Ala	GCT	20.23	0.28	Thr	ACT	13.24	0.23	Ser	TCT	13.51	0.18	Pro	CCT	18.03	0.29
Ala	GCC	28.43	0.40	Thr	ACC	21.52	0.38	Ser	TCC	17.37	0.23	Pro	CCC	20.51	0.33

Analýza využití kodonů (codon usage)

Gly	GGG	17.08	0.23	Arg	AGG	12.09	0.22	Trp	TGG	14.74	1.00	Arg	CGG	10.40	0.19
Gly	GGA	19.31	0.26	Arg	AGA	11.73	0.21	End	TGA	2.64	0.61	Arg	CGA	5.63	0.10
Gly	GGT	13.66	0.18	Ser	AGT	10.18	0.14	Cys	TGT	9.99	0.42	Arg	CGT	5.16	0.09
Gly	GGC	24.94	0.33	Ser	AGC	18.54	0.25	Cys	TGC	13.86	0.58	Arg	CGC	10.82	0.19
Glu	GAG	38.82	0.59	Lys	AAG	33.79	0.60	End	TAG	0.73	0.17	Gln	CAG	32.95	0.73
Glu	GAA	27.51	0.41	Lys	AAA	22.32	0.40	End	TAA	0.95	0.22	Gln	CAA	11.94	0.27
Asp	GAT	21.45	0.44	Asn	AAT	16.43	0.44	Tyr	TAT	11.80	0.42	His	CAT	9.56	0.41
Asp	GAC	27.06	0.56	Asn	AAC	21.30	0.56	Tyr	TAC	16.48	0.58	His	CAC	14.00	0.59
Val	GTG	28.60	0.48	Met	ATG	21.86	1.00	Leu	TTG	11.43	0.12	Leu	CTG	39.93	0.43
Val	GTA	6.09	0.10	Ile	ATA	6.05	0.14	Leu	TTA	5.55	0.06	Leu	CTA	6.42	0.07
Val	GTT	10.30	0.17	Ile	ATT	15.03	0.35	Phe	TTT	15.36	0.43	Leu	CTT	11.24	0.12
Val	GTC	15.01	0.25	Ile	ATC	22.47	0.52	Phe	TTC	20.72	0.57	Leu	CTC	19.14	0.20
Ala	GCG	7.27	0.10	Thr	ACG	6.80	0.12	Ser	TCG	4.38	0.06	Pro	CCG	7.02	0.11
Ala	GCA	15.50	0.22	Thr	ACA	15.04	0.27	Ser	TCA	10.96	0.15	Pro	CCA	17.11	0.27
Ala	GCT	20.23	0.28	Thr	ACT	13.24	0.23	Ser	TCT	13.51	0.18	Pro	CCT	18.03	0.29
Ala	GCC	28.43	0.40	Thr	ACC	21.52	0.38	Ser	TCC	17.37	0.23	Pro	CCC	20.51	0.33

Navrhování sekvencí primerů pro PCR

- Standardní primery
- Modifikované oligonukleotidy na 5'-konci pro klonování
- Oligonukleotidy jako hybridizační sondy pro real-time PCR
 - ◆ specifčnost
 - ◆ jedinečnost





PCR - Syntéza obou řetězců u specifické sekvence



Výběr vhodné strategie před návrhem primerů

- K čemu jsou primery určeny
 - ◆ Standardní end-point PCR
 - ◆ Degenerovaná PCR
 - ◆ Multiplex PCR
 - ◆ Sekvenování (primer walking)
 - ◆ Real-time PCR
 - ◆ Detekce jednonukleotidových polymorfizmů (SNP) nebo variací
 - ◆ Studium metylace
 - ◆ Sondy pro microarray
- Z jakých dat vycházíme
 - ◆ Jednoduchá sekvence DNA / proteinu
 - ◆ Sekvenční příložená DNA / proteinu
 - ◆ GenBank ID/Gene ID/rsSNP ID
 - ◆ Optimální je využívat kompletní/co nejdelší templát

Pravidla pro design primeru pro PCR

- Relativně snadná výpočetní záležitost –
prohledávání sekvence a identifikace krátkých
sekvencí splňujících určitá kritéria
 - ◆ Délka primeru
 - ◆ Obsah G+C
 - ◆ Teplota T_m
 - ◆ Specificita
 - ◆ Komplementarita primerových sekvencí
 - ◆ Sekvence 3'-konce

Jedinečnost primeru



- Na jedinečnost primeru a jeho hybridizační vlastnosti (annealing) má vliv délka primeru a velikost templátové DNA
 - ◆ Délka (17 – 28 bází dlouhé)
- Možná hybridizační místa primeru by se také neměla nacházet na DNA tvořících případné kontaminace vzorků

Templátová DNA

5' . . . TCAACTTAGCATGATCGGGTA . . . GTAGCAGTTGACTGTACAACCTCAGCAA . . . 3'
TGCTAAGTTG CAGTCAACTGCTAC TCGT AGTTG
A

Primer 1 5' – TGCTAAGTTG – 3'

Není jedinečný!

Primer 2 5' – CAGTCAACTGCTAC – 3'

Jedinečný!

Zastoupení bází

- Zastoupení bází ovlivňuje vlastnosti hybridizace a reasociace primeru
- Žádoucí je náhodná distribuce bází bez oblastí bohatých na AT nebo GC
- Obvyklý obsah G+C, který poskytuje stabilní hybridy je 40-60 %, ale závisí také na obsahu G+C templátu

Templátová DNA

5' ...TCAACTTAGCATGATCGGGCA...AAGATGCACGGGCCTGTACACAA...3'

TGGCCTAGCATGCT TGGCCTAGCATGCT

Teplota T_m (Melting temperature)

- ◆ mají T_m teplotu 50 – 65 °C

$$T_a = 0,3 \times T_m^{\text{Primer}} + 0,7 \times T_m^{\text{Produkt}} - 25$$

kde T_m^{Primer} je hodnota T_m nejméně stabilního páru primer-matrice a T_m^{Produkt} je hodnota T_m amplifikačního produktu.

- Orientačně lze vypočítat T_a podle vztahu:

$$T_m = 2(A+T) + 4(G+C)$$

$$T_a = T_m - 5 \text{ °C}$$

Vnitřní sekvence a struktura primeru

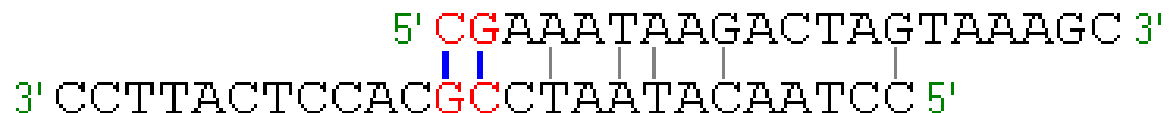
- Stabilita duplexů oligonukleotidů a vlásenek se uvádí v ΔG (kcal/mol)
- Oligonukleotidy nejsou komplementární navzájem na 3'-koncích, takže nevytvářejí navzájem nebo samy se sebou duplexy
- Netvoří vnitřní sekundární struktury



- ◆ Chybně navržená dvojice primerů, která vytváří stabilní duplex na 3'-konci:



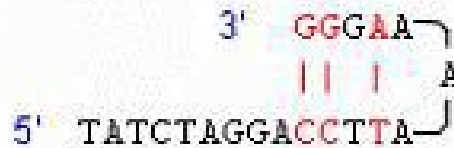
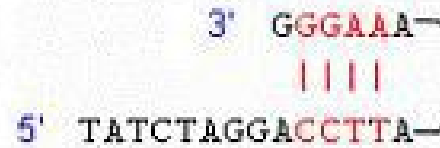
- ◆ Správně navržená dvojice primerů, která vytváří pouze málo stabilní duplex na 5'-konci; na 3'-konci je G nebo C zaručující stabilní párování s templátem:



- ◆ Chybně navržený primer, vytvářející vlásenku:

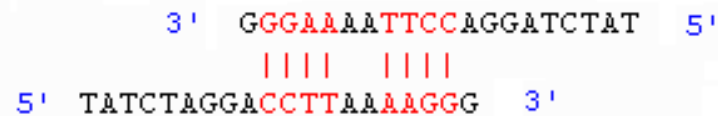


Hairpin

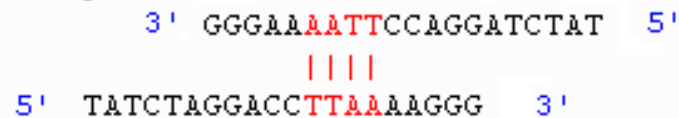


Self-Dimer

8 bp

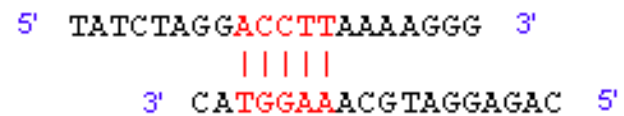


4 bp



Dimer

forward primer




reverse primer

GC svorky a 3' - koncová stabilita

■ GC svorka

- ◆ Přítomnost G nebo C mezi posledním 4 bázemi na 3'-konci primeru
- ◆ Zásadní pro zvýšení prevence falešného prodlužování a zvýšení specifičnosti primeru
- ◆ >3 G nebo C v blízkosti 3'-konce jsou však nežádoucí

5' GAAGTACGGAAGAAGC 3'
CTTTAAACCCTTCATGCCTTCTTCGACACCTAAATGGTCTAATTTTCAGCTCC



Jedinečnost primerů

na matricové DNA nemají falešná vazebná místa

- Nesprávně navržený primer s falešnými vazebnými místy na templátové DNA:

```

5'(1029) AAGGCTAGAGAAAAATATGG (1048)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(948) tttcttacccttttt-tacc (966)5'
    
```

```

5'(1029) AAGGCTAGAGAAAAATATGG (1048)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(1191) tttgtattgcattatatacc (1210)5'
    
```

```

5'(1029) AAGGCTAGAGAAAAATATGG (1048)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(395) tccatttttcttttttatctt (414)5'
    
```

- Správně navržený primer, který nemá falešná vazebná místa na templátu:

```

5'(2476) CCTAACATAATCCGCACCTCATTC (2452)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(787) taaatctattagttacacataacc (811)5'
    
```

```

5'(2476) CCTAACATAATCCGCACCTCATTC (2452)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(3211) caattgtaactataactgcggtatc (3235)5'
    
```

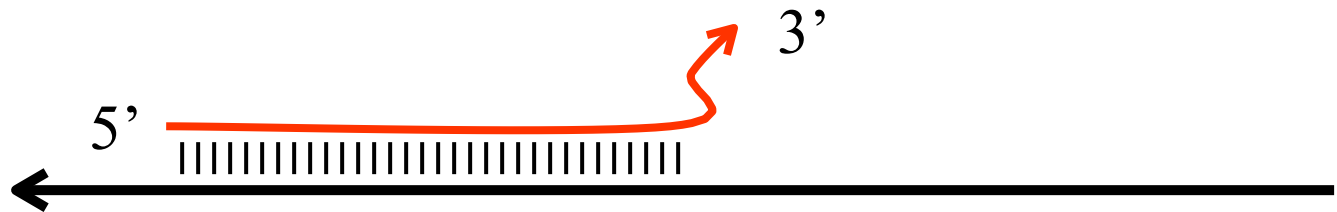
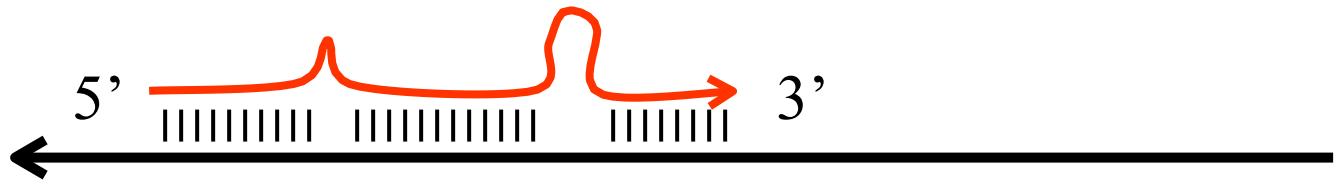
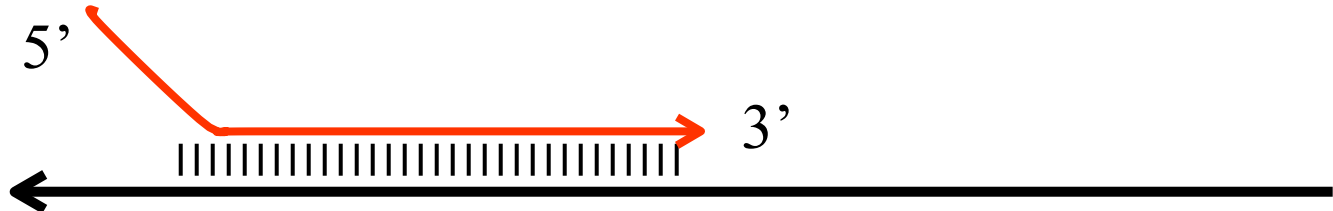
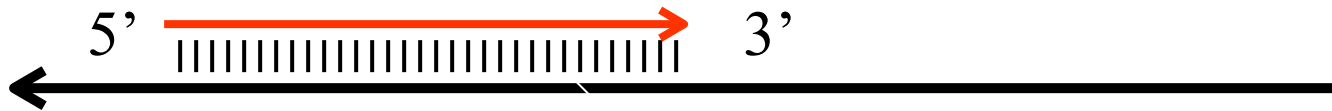
```

5'(2476) CCTAACATAATCCGCACCTCATTC (2452)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(1194) gtattgcattatataacctctgtag (1218)5'
    
```

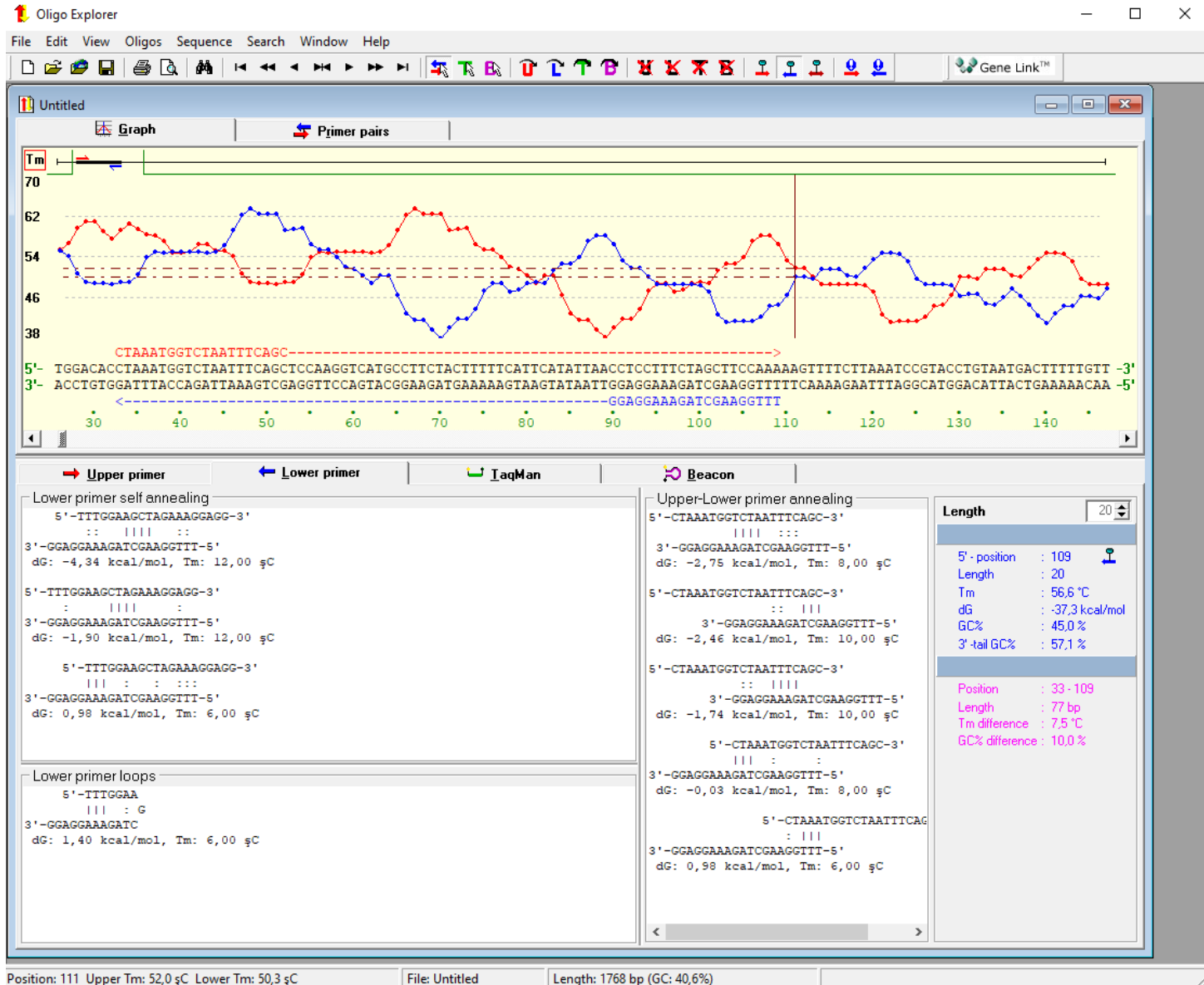
```

5'(2476) CCTAACATAATCCGCACCTCATTC (2452)3'
         | | | | | | | | | | | | | | | | | | | | | |
3'(1469) atattgta-tatacgaactaaatct (1492)5'
    
```

Kdy je primer ještě primerem?



Pro návrh primerů se obvykle používá specializovaný software



Počítačový návrh primerů

- Umožňuje řada molekulárně biologických programů
- Některé jsou volně dostupné na internetu
 - ◆ Primer3 (<http://primer3.sourceforge.net/webif.php>)
 - ◆ Primer3Plus (<http://www.bioinformatics.nl/cgi-bin/primer3plus/primer3plus.cgi>)
 - ◆ PrimerZ – design oligo pro lidské promotory, exony a SNP (<http://genepipe.ngc.sinica.edu.tw/primerz/beginDesign.do>)
 - ◆ PerlPrimer
 - ◆ BioTools
 - ◆ WebPrimer
- Kalkulátory vlastností primerů
 - ◆ IDT Oligo Analyzer (<http://eu.idtdna.com/SciTools/SciTools.aspx?cat=DesignAnalyze>)
 - ◆ BioMath (<http://www.promega.com/biomath/calc11.htm>)
 - ◆ PrimerBlast
 - ◆ UCSC In-Silico PCR
 - ◆ AutoDimer

Oligo Calculator

Oligo Calc: Oligonucleotide Properties Calculator

Enter Oligonucleotide Sequence Below
OD calculations are for single-stranded DNA or RNA

[Nucleotide base codes](#)

Reverse Complement Strand(5' to 3') is:

5' [modification](#) (if any) 3' [modification](#) (if any) Select molecule

nM Primer Measured Absorbance at 260 nanometers

mM Salt (Na⁺)

CalculateSwap StrandsBLASTmfold

Physical Constants

Length: Molecular Weight: GC content: %

1 ml of a sol'n with an Absorbance of at 260 nm
is microMolar ⁵ and contains micrograms.

Melting Temperature (T_M) Calculations

°C (Basic)

°C (Salt Adjusted)

°C (Nearest Neighbor)

Thermodynamic Constants Conditions: 1 M NaCl at 25°C at pH 7.

RlnK cal/(°K*mol)

deltaH Kcal/mol

deltaG Kcal/mol

deltaS cal/(°K*mol)

Deprecated Hairpin/self dimerization calculations

(Minimum base pairs required for single primer self-dimerization)

(Minimum base pairs required for a hairpin)

Check Self-Complementarity

Primer 3 <http://primer3.sourceforge.net/webif.php>

Primer3 Input (version 0.4.0) - Mozilla Firefox

Soubor Úpravy Zobrazení Historie Záložky Nástroje Nápověda

http://frodo.wi.mit.edu/primer3/input.htm

Primer3 Input (version 0.4.0)

Primer3

(v. 0.4.0) Pick primers from a DNA sequence.

[Checks for mispriming in template.](#) [disclaimer](#) [Primer3 Home](#)
[Primer3plus interface](#) [cautions](#) [FAQ/WIKI](#)

Paste source sequence below (5'→3', string of ACGTNacgtn -- other letters treated as N -- numbers and blanks ignored). FASTA format ok. Please N-out undesirable sequence (vector, ALUs, LINEs, etc.) or use a [Mispriming Library \(repeat library\)](#):

```
>SA44kb001 [org=Staphylococcus aureus] [strain=CCM 885] [clone=7/IV] Staphylococcus aureuss
EcoRI-clone from common 44 kb SmaI fragment
GAATTCAAAACCCAGCAAAAAGCTGTGAAAAAGCCATTACCAAGTAAAGATAATTTGGCTATATTGTATGGAGAAGGATTTTCATATTTGTAAGGCG
AATTATTTGGAAAACATCGACATGGTGAAGATTGTCTGTTCTGTTTAAAGGTTTTAAGTGATTAATCAAGCACACTCAAATAGTGTTATAATTAT
AAATGAATATGGTTTGGATAAGTCTGAGACAATGCATGTTTCAGGCTTTAATTGTGTATAAAAGTTTTGGTGATTGCATAAGAGATGGCGGTAATA
AATGTTATTATTAAGTGTGCACGCAGTATCATTAGTTATAAAAATGTAGCTGTTAAAAAGTCAAAAATACATCGAATGTAGTTAGGCATATAATATA
```

<input checked="" type="checkbox"/> Pick left primer, or use left primer below:	<input type="checkbox"/> Pick hybridization probe (internal oligo), or use oligo below:	<input checked="" type="checkbox"/> Pick right primer, or use right primer below (5' to 3' on opposite strand):
<input type="text"/>	<input type="text"/>	<input type="text"/>

[Sequence Id:](#) A string to identify your output.

[Targets:](#) E.g. 50,2 requires primers to surround the 2 bases at positions 50 and 51. Or mark the [source sequence](#) with [and]: e.g. ...ATCT[CCCC]TCAT.. means that primers must flank the central CCCC.

Hotovo

Pick Primers Reset Form

Sequence Id: A string to identify your output.

Targets: E.g. 50,2 requires primers to surround the 2 bases at positions 50 and 51. Or mark the source sequence with [and]: e.g. ...ATCT[CCCC]TCAT.. means that primers must flank the central CCCC.

Excluded Regions: E.g. 401,7 68,3 forbids selection of primers in the 7 bases starting at 401 and the 3 bases at 68. Or mark the source sequence with < and >: e.g. ...ATCT<CCCC>TCAT.. forbids primers in the central CCCC.

Product Size Ranges: 150-250 100-300 301-400 401-500 501-600 601-700 701-850 851-1000

Number To Return: Max 3' Stability:

Max Repeat Mispriming: Pair Max Repeat Mispriming:

Max Template Mispriming: Pair Max Template Mispriming:

Pick Primers Reset Form

General Primer Picking Conditions

Primer Size Min: Opt: Max:

Primer Tm Min: Opt: Max: Max Tm Difference: Table of thermodynamic parameters:

Product Tm Min: Opt: Max:

Primer GC% Min: Opt: Max:

Oligo – příklad komerčního software

Oligo 7 Demo - Human eIF-4E.seq

File Edit Analyze Search Select Change View Window Help

Sequence

File: Human eIF-4E.seq

DNA Sequence		Selected Oligo	Position	Length
Sequence Length:	1868 nt	<input checked="" type="checkbox"/> Forward Primer	997	22
Reading Frame:	+1	<input checked="" type="checkbox"/> Reverse Primer	1061	21
Current Oligo Length:	21 nt	<input checked="" type="checkbox"/> Upper Oligo	956	21
Position:	956	<input type="checkbox"/> Lower Oligo	---	---
t_m :	49.1°C	<input checked="" type="checkbox"/> PCR Product	[85,---] nt	

#	Feature	Location
1	source	-18..1850

pos: tm:

950 960 970 980 990 1000 1010 1020 1030 1040 1050 1060 1070 1080

ACATACAGATTTTACCTATCC · TGGCATTCTATACTTTACAGG ·

ATTACCATTAATTACATACAGATTTTACCTATCCACAATAGTCAGAAAAACAATTGGCATTCTATACTTTACAGGAAAAAAAAATTCTGTTGTTCCATTTTATGCAGAAGCATATTTTGCTGGTTTGAAAAGATTATGATGCAT
 TAATGGTAATTAATGTATGTCTAAAATGGATAGGTGTTATCAGTCTTTTGTGAACCGTAAAGATATGAAATGTCTTTTTTTAAGACAACAAGGTAAAATACGTCTTCGTATAAAAACGACCAAACCTTTCTAATACTACGTA

CGACCAAACCTTTCTAATACTA

I T I N Y I Q I L P I H N S Q K T T W H F Y T L Q E K K F C C S I L C R S I F C W F E R L - C I

Ready...

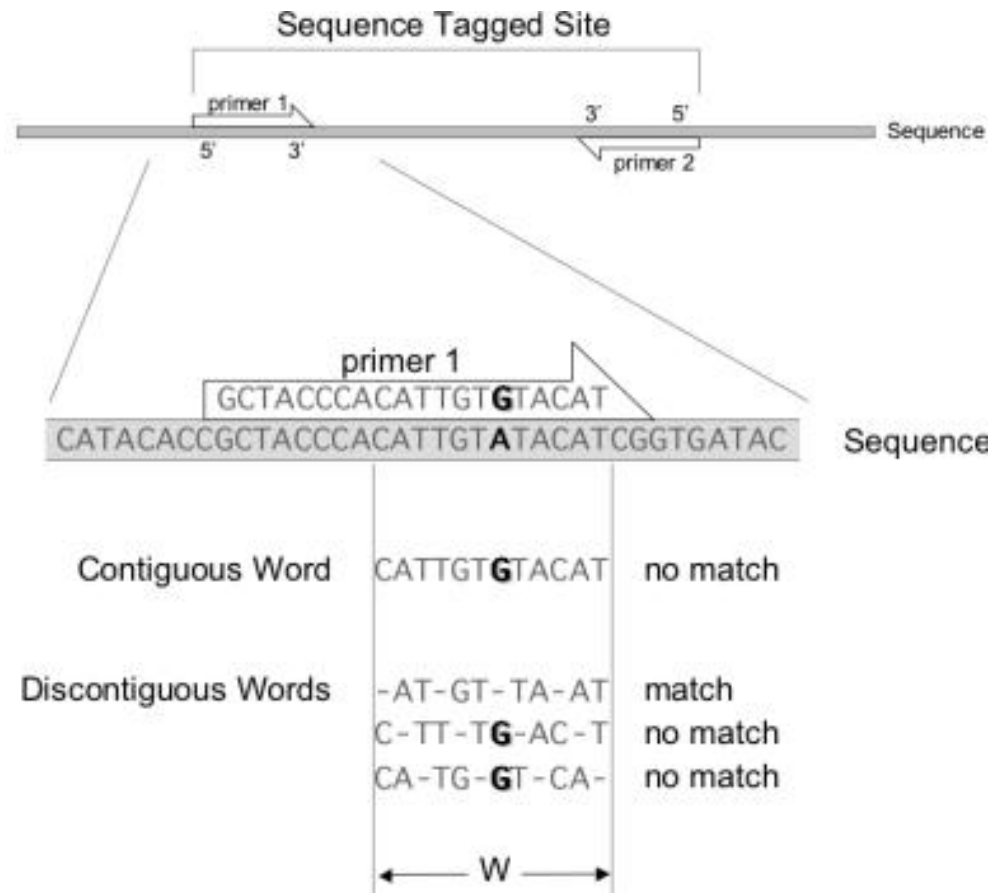
Výsledky poskytované programy pro návrh oligonukleotidů

- Výběr optimálního páru primerů
- Sekvence primerů
- Délka primerů a hodnota T_m
- Velikost produktu
- Posouzení sekundárních struktur
- Podmínky reakce
- Alternativní primery

Electronic PCR (e-PCR)



- Většina programů pro design PCR-primerů je omezená délkou templátu
- Posouzení falešných vazebných míst ideálně vyžaduje práci s celogenomovou sekvencí nebo částí databáze
- e-PCR označuje výpočetní postup, který se používá k prohledávání sekvencí DNA s cílem nalezení jedinečných míst, které odpovídají výsledkům PCR
- Nástroje lze využít pro eliminaci falešných pozitivních výsledků nebo design degenerovaných univerzálních oligonukleotidů



PCR Primer Mapping – UCSC In-Silico PCR

<http://genome.ucsc.edu/cgi-bin/hgPcr?db=mm9>

Home Genomes Blat Tables Gene Sorter Session FAQ Help

UCSC In-Silico PCR

Genome: Assembly: Forward Primer: Reverse Primer:

Max Product Size: Min Perfect Match: Min Good Match: Flip Reverse Primer:

About In-Silico PCR

In-Silico PCR searches a sequence database with a pair of PCR primers, using an indexing strategy for fast performance.

Configuration Options

Genome and Assembly - The sequence database to search.

Forward Primer - Must be at least 15 bases in length.

Reverse Primer - On the opposite strand from the forward primer. Minimum length of 15 bases.

Max Product Size - Maximum size of amplified region.

Min Perfect Match - Number of bases that match exactly on 3' end of primers. Minimum match size is 15.

Min Good Match - Number of bases on 3' end of primers where at least 2 out of 3 bases match.

Flip Reverse Primer - Invert the sequence order of the reverse primer and complement it.

Output

When successful, the search returns a sequence output file in fasta format containing all sequence in the database that lie between and include the primer pair. The fasta header describes the region in the database and the primers. The fasta body is capitalized in areas where the primer sequence matches the database sequence and in lower-case elsewhere. Here is an example:

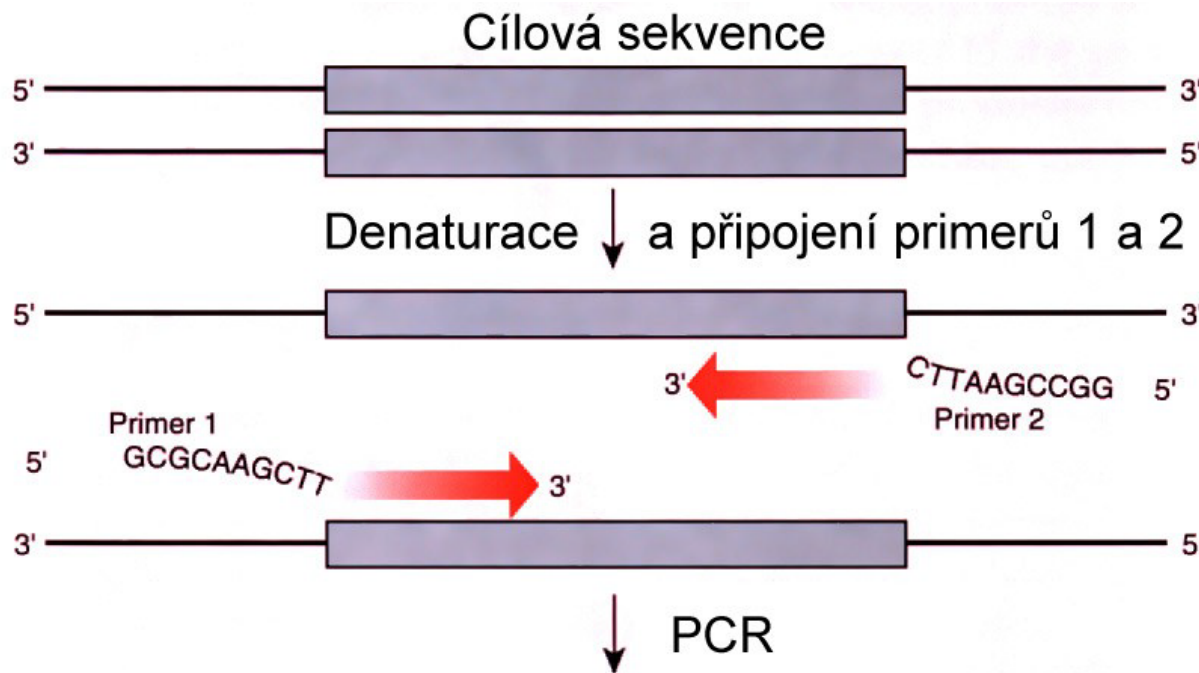
```
>chr22:31000551+31001000 TAACAGATTGATGATGCATGAAATGGG CCCATGAGTGGCTCCTAAAGCAGCTGC
TtACAGATTGATGATGCATGAAATGGGgggtggccaggggtgggggtga
gactgcagagaaaggcagggtggttcataacaagctttgtgogtcccaa
tatgacagctgaagttttccaggggtgatggtgagccagtgagggtaa
tacacagaacatcctagagaaacccctattccttaagattaaaaataaa
```


Pokročilý návrh primerů

- Alelově specifické primery
- Molekulární diagnostika
 - ◆ Vícenásobné detekce - primery pro multiplex PCR
 - ◆ Zajištění kompatibility primerů v reakci
- Konsenzní primery
 - ◆ Vyžaduje identifikaci konzervativních oblastí na základě mnohonásobných přiložení sekvencí (multiple alignment)
 - ◆ Pro klonování
 - ◆ Pro mutagenezi
- Primery pro modifikaci konců produktů PCR



Modifikace konců DNA, Připojení sekvencí prostřednictvím 5'-konců primerů



„sticky foot“



■ Přidávané sekvence

- ◆ RE místa
- ◆ Promotory
- ◆ Terminátory
- ◆ Translační signály

Zdroje pro návrh multiplex PCR

- NCBI/ Primer-BLAST
- MultiPLX (<http://bioinfo.ebc.ee/multiplx/>)
- PrimerStation (<http://ps.cb.k.u-tokyo.ac.jp/index.html>)
 - ◆ Lidský genom
 - ◆ Specifikace exonů
 - ◆ Vyloučení variabilních oblastí se SNP
- Oligo Explorer (<http://www.genelink.com/tools/gl-oe.asp>)
 - ◆ Posouzení dimerů primerů v multiplexovém uspořádání

Webové zdroje pro design primerů pro real-time PCR

- NCBI Probe Database
- RTPrimerDB
- Primer Bank
- qPrimerDepot
- PCR-QPPD
- PerlPrimer
- Komerční databáze (např. ROCHE,...)

Nejčastěji používané softwarové balíky pro manipulaci se sekvencemi

- Geneious (Biomatters, Inc., New Zealand)
- Ugene (<http://ugene.net/>) freeware
- CLC Genomics Workbench (CLC bio, Cambridge)
- Vector NTI[®] (Life Technologies, Carlsbad, CA)
- Bioinformatics Toolbox rozšíření pro MATLAB[®]
- Hitachi DNASIS[®] MAX Sequence Analysis Software (Helixx Technologies, Inc., Canada)
- DNASTAR Lasergene (DNASTAR, Inc., Madison, WI)
- Accelrys GCG Package (Accelrys Inc., San Diego)