



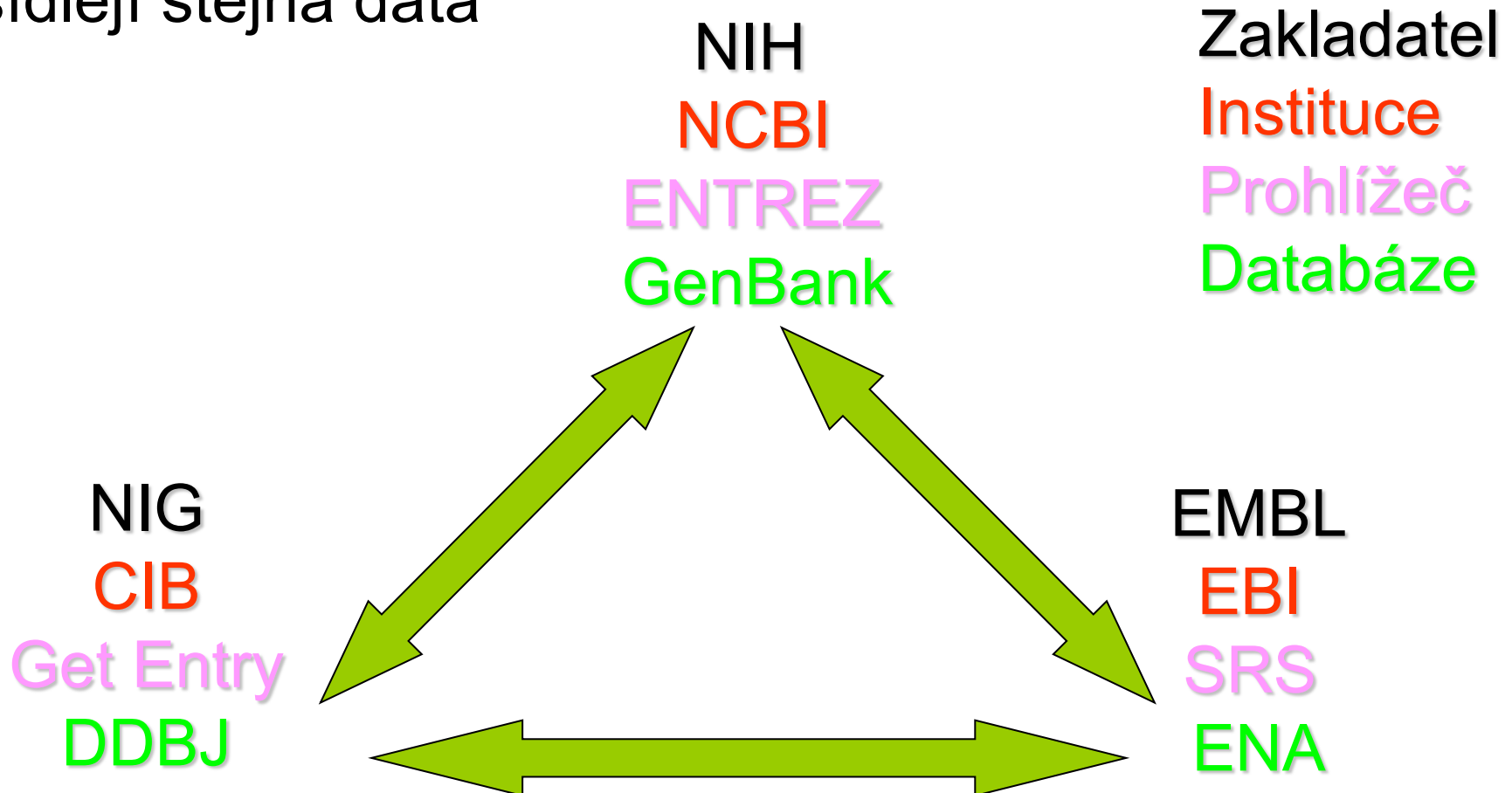
Bi5000 – Bioinformatika

Zpracování sekvenčních dat pro zaslání do nukleotidových databází



Mezinárodní spolupráce primárních sekvenčních databází

Všechny tři hlavní bioinformatická centra (databáze)
sídlí stejná data



Divize GenBank



<https://www.ncbi.nlm.nih.gov/genbank/htgs/divisions/>

<ftp://ftp.ncbi.nlm.nih.gov/genbank>

GenBank Database Divisions

GenBank divisions are divided into two general categories and were described in an (Genome Research (1997) 7(10)) article by Ouellette and Boguski; the full-text article is available ([Database Divisions and Homology Search Files: A Guide for the Perplexed](#)). The "Organismal" category includes databases pertaining to sequences derived from specific organisms and the "Functional" databases pertain to different types of sequence data being collected. Sequence records exist only in one GenBank division. For example, the HTG division includes unfinished sequences (phases 0, 1, and 2) being generated from several different organisms. As a sequence is updated to phase 3, it is moved into the appropriate organismal division. For instance, human phase 3 (finished) HTG sequences are located in the PRI division. The GenBank divisions listed here represent the location of the annotated sequence records; for homology search purposes the records are reformatted and stored in the [BLAST databases](#). The different database divisions currently available, as well as the related BLAST database, are listed below. An example of a submission (one accession number) that has progressed through phase 1, phase 2, and phase 3 is available ([Examples](#)).

HTGs Resource:

- [About HTGs](#)
- [Submitting HTGs](#)
- [Processing HTGs](#)
- [HTGs FAQ](#)

Organismal Divisions:

Database	Division	BLAST	Example
BCT	Bacterial sequences	nr, month	
PRI	Primate sequences	nr, month	Human Phase 3
ROD	Rodent sequences	nr, month	
MAM	Other mammalian sequences	nr, month	
VRT	Other vertebrate sequences	nr, month	
INV	Invertebrate sequences	nr, month	Drosophila, C. elegans Phase 3
PLN	Plant and Fungal sequences	nr, month	Arabidopsis Phase 3
VRL	Viral sequences	nr, month	
PHG	Phage sequences	nr, month	
RNA	Structural RNA sequences	nr, month	
SYN	Synthetic and chimeric sequences	nr, month	
UNA	Unannotated sequences	nr, month	

← Divize podle organizmů

Functional Divisions:

Database	Division	BLAST	Example
EST	Expressed Sequence Tags	dbest, month	
STS	Sequence Tagged Sites	dbsts, month	
GSS	Genome Survey Sequences	dbgss, month	
HTG	High Throughput Genomic sequences	htgs, month	All Organisms: Phase 0, 1, and 2

← Funkční divize

Postup zpracování sekvenčních dat



- **Primární analýza dat ze sekvenátoru**

- Zachycení obrazu
- Zpracování obrazu
- Base calling
- Provedení kontroly kvality

Sequence Read Archive:
Nezpracovaná sekvenační
data pro zvýšení
reprodukovatelnosti a
usnadnění nových objevů.

- **Sekundární analýza dat**

- Alingment (RefSeq data)
 - Mismatches, validace,
vizualizace
- Assembly -> kontigy/scaffoldy

Whole Genome Shotgun:
sestavy neúplných nebo
kompletních genomů a
chromozomů prokaryot a
eukaryot získané NGS.

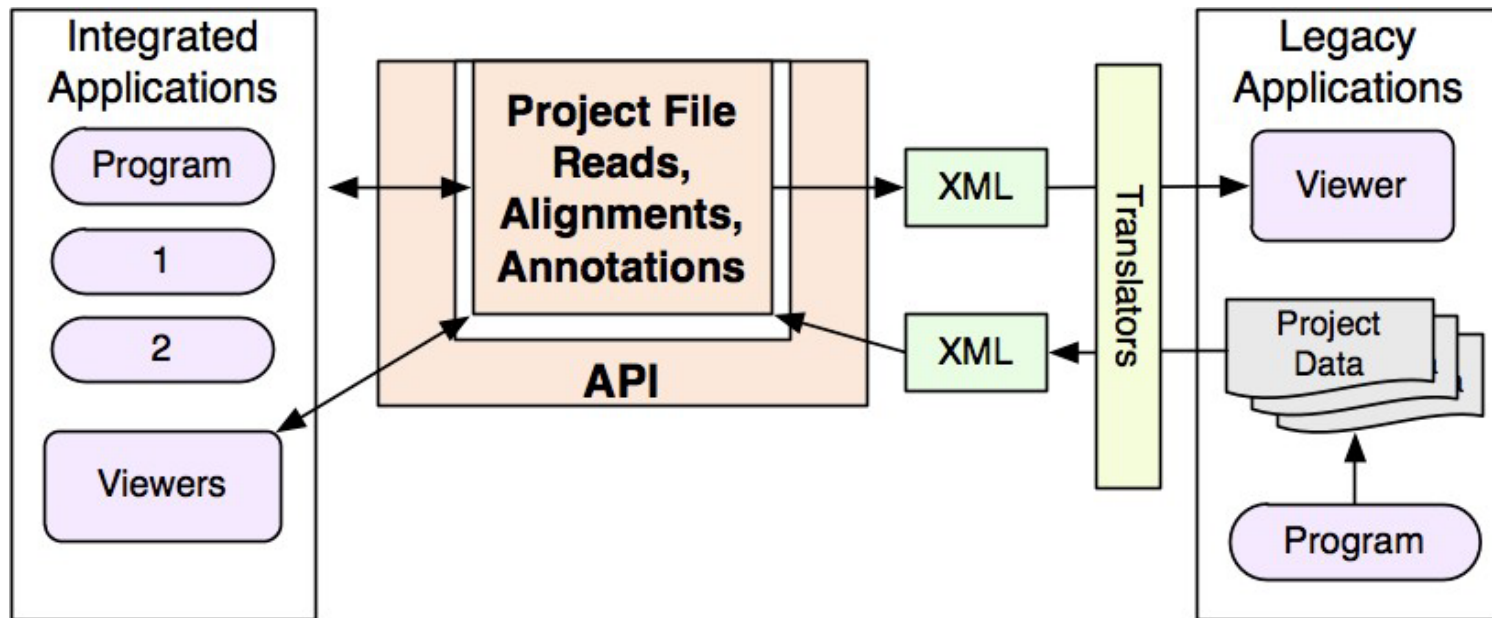
- **Terciární analýza - automatická anotace**

- Anotace, Geny, Variace,
Diferenciální exprese, Metylace, atd.

Řada specializovaných
databází: GenBank,
Genomy, SNP, GEO

Komplexní přístup vs. manuální zpracování

Komplexní přístupy vyžadují strukturovaná, indexovaná a anotovaná data



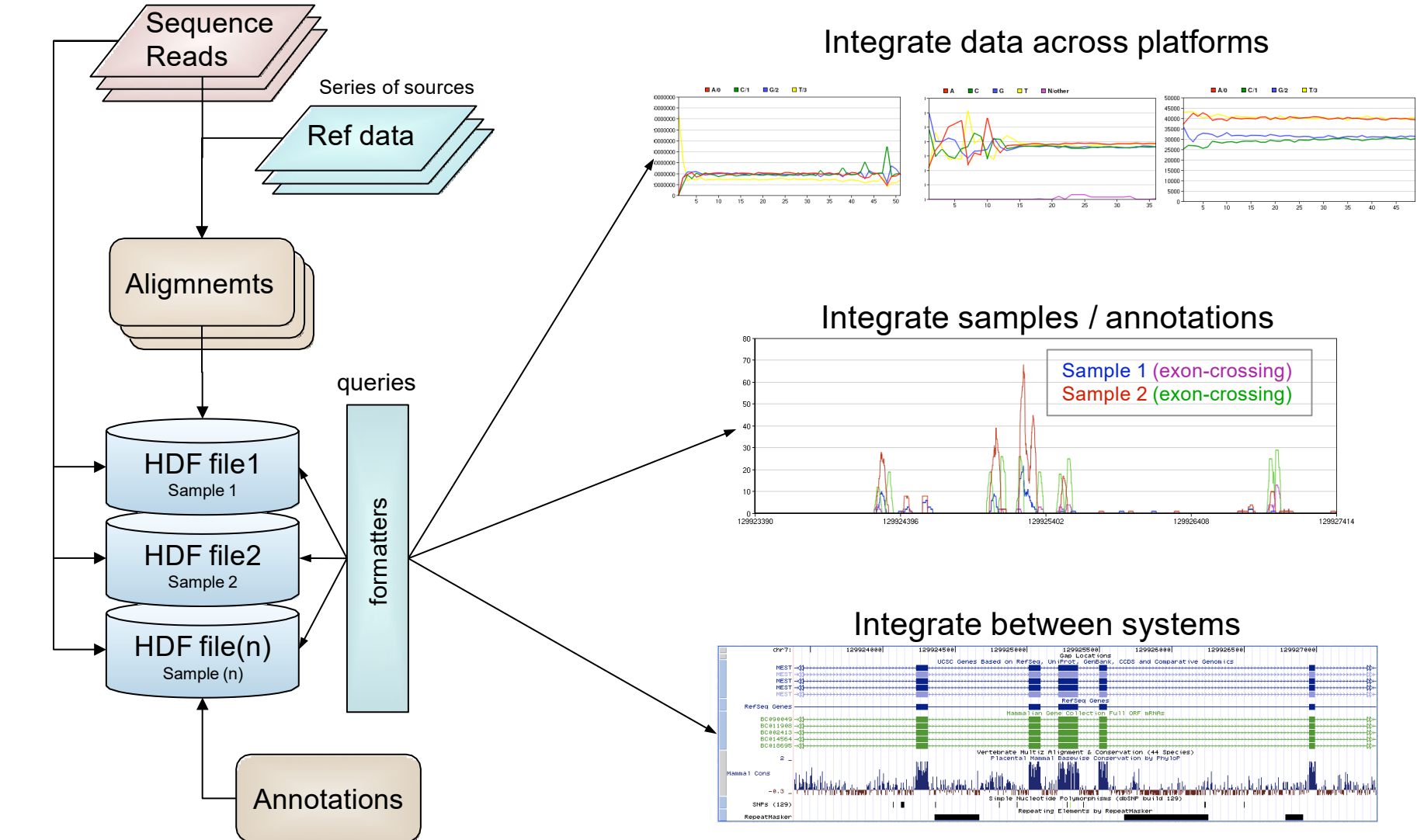
- Zpracování velkého množství dat
- Integrace více typů dat
- Porovnávání vzorků mezi sebou

Formáty sekvenačních dat z NGS



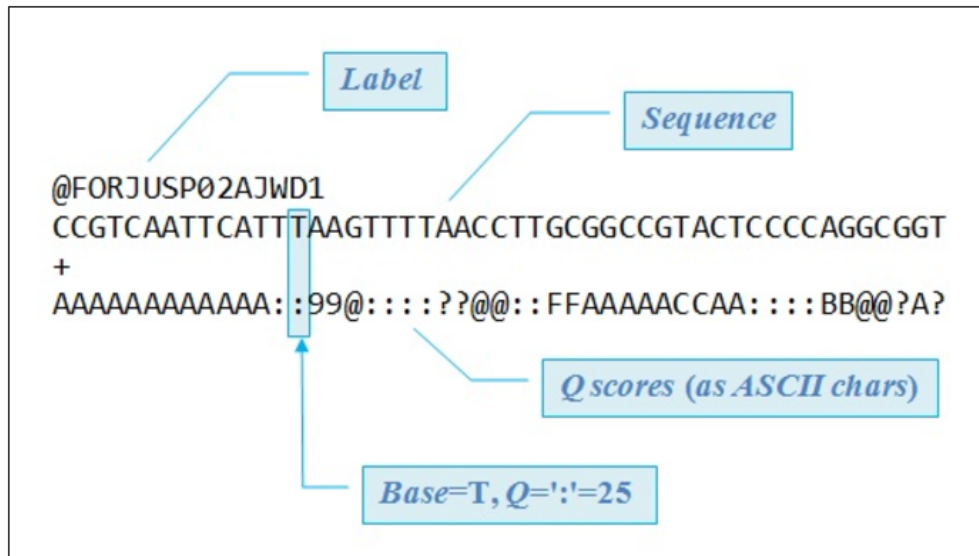
- FASTQ
 - Illumina
 - Standardní formát akceptovaný databázemi
- SFF, BAM, qual
 - starší formáty používané 454 a Ion Torrent
 - Obsahují alignment jednotlivých čtení
 - Většinou lze konvertovat do FASTQ
- fast5, HDF5 - Hierarchické formáty dat
 - formáty pro dlouhá čtení
 - PacBio, Nanopore

Hlubší integrace dat do komplexních databází (BioHDF)



FASTQ formát

- Univerzální formát akceptovaný databázemi
- Každý záznam obsahuje 4 řádky
 - řádek 1 začíná hlavičkou '@'ID sekvence
 - Řádek 2 obsahuje primární sekvenci
 - Řádek 3 začíná '+' a může následovat stejné ID a popis
 - Řádek 4 obsahuje zakódované hodnoty o kvalitě sekvence a musí obsahovat stejný počet znaků jako řádek 2



Struktura sekce Illumina Header

```
@HWUSI-  
EAS611:34:6669YAAXX:5:1:5069:1159  
1:N:0:
```

- Starts with @ (required by fastq spec)
- Instrument ID (HWUSI-EAS611)
- Run number (34)
- Flowcell ID (6669YAAXX)
- Lane (5)
- Tile (1)
- X-position (5069)
- Y-position (1159)
- [space]
- Read number (1)
- Was filtered (Y/N) (N) - You wouldn't normally see the Ys
- Control number (0 = no control)
- Sample number (only if demultiplexed using Illumina's software)



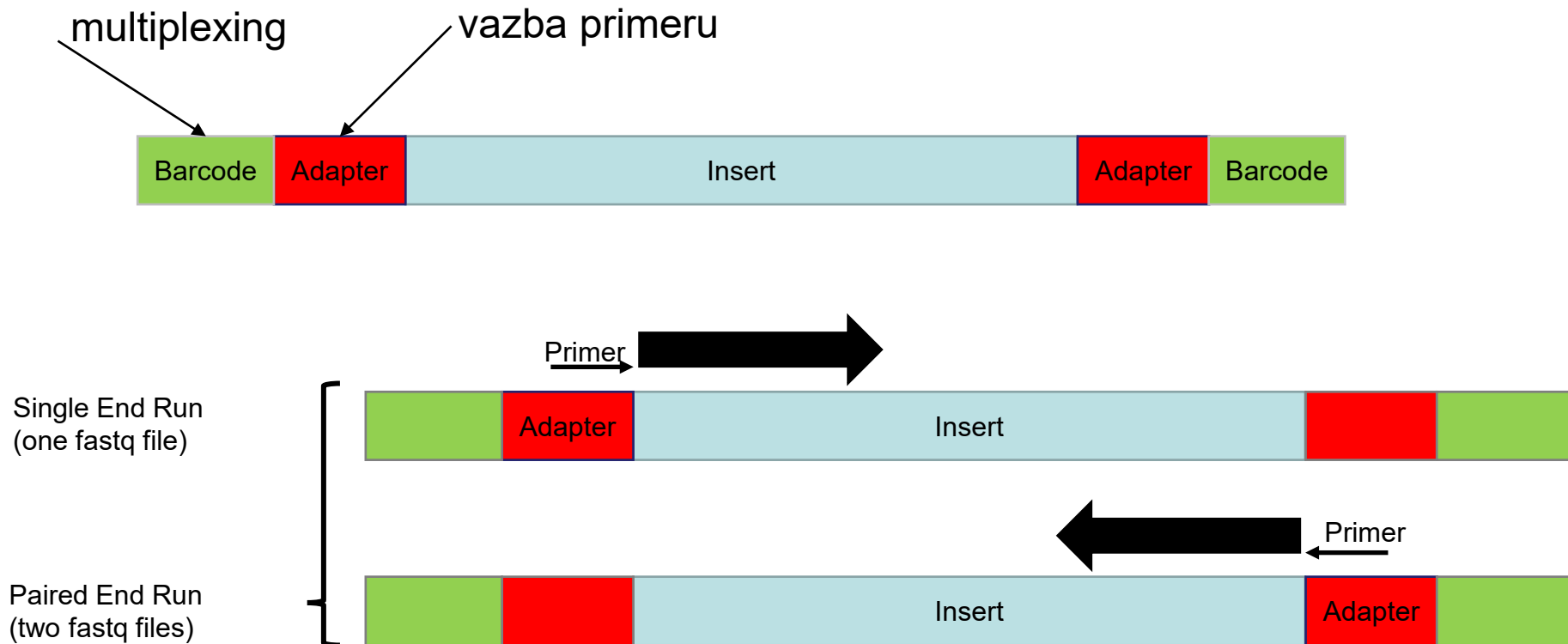
HDF5 a Fast5 formáty

- PacBio a Nanopore
- HDF (hierarchical data format)
 - zkomprimovaná složka s hierarchickou strukturou
 - Skupiny
 - Datové sady
 - Atributy
- Fast5 - komplexní formát
 - odvozený z HDF5
 - má definovanou konkrétní hierarchii souborů
- Obsahuje data pro base calling
 - vyšší informační obsah, base calling může být provedeno opakovaně
 - pro složení sekvence může být převeden do FASTQ

The screenshot shows the HDFView application window. The file path is `/Users/msw/Desktop/seqc_brain_3.txt.h5`. The left sidebar shows a tree view with folders like `alignments`, `genome`, `sequences`, and `seqid`. The main window displays a table with columns `ref_id`, `beg_pos`, `end_pos`, and `num_read`. The table contains 22 rows of data. Below the table, there is a status bar showing `cluster (1136145282)` and `Compound/Vdata, 1270853`.

	ref_id	beg_pos	end_pos	num_read
0	22	5705	14703	454656
1	22	661	5649	246223
2	0	556117	559948	244852
3	0	554323	556061	63856
4	22	14712	16570	47040
5	4	134290167	134291100	41455
6	4	79982071	79983331	32404
7	16	48538223	48538748	27540
8	2	97819217	97819732	20338
9	2	97818721	97819082	18158
10	4	134291101	134291462	14506
11	1	49310335	49310542	14429
12	5	62341988	62342229	14024
13	10	10486225	10486552	12565
14	10	10486762	10487151	12228
15	23	125433390	125433757	11051
16	1	87905671	87905992	10865
17	4	134288734	134289496	10687
18	0	559952	560170	9975
19	4	79981593	79981750	9813
20	17	43633622	43633808	9797
21	6	45258094	45258249	9147

Jak čteme sekvence z NGS?



1. Demultiplexing
2. Base Call Quality
3. Adapter Content (trimming)
4. Mapping Quality



QC Metrics



Kontrola kvality

- Statistiky čtení
 - počet čtení a medián jejich délky
 - obsah G+C
 - skóre kvality bází a jeho distribuce
- Kontaminace
 - adaptory, primery
 - hostitelská DNA u nebuněčných genomů
 - křížové kontaminace
- Hloubka pokrytí
 - Volba vyžadované hloubky závisí na aplikaci
 - pro celogenom je potřeba minimálně 30-50x pokrytí u krátkých čtení a minimálně 100-300x u dlouhých čtení
- FastQC – příklad nástroje pro hodnocení kvality u Illumina
<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Základní statistiky

- Illumina

Filename	CCM9024_R1.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	2761796
Sequences flagged as poor quality	0
Sequence length	150
%GC	34

Hodnotí se především:

- počet čtení
- délka čtení
- distribuce délek čtení

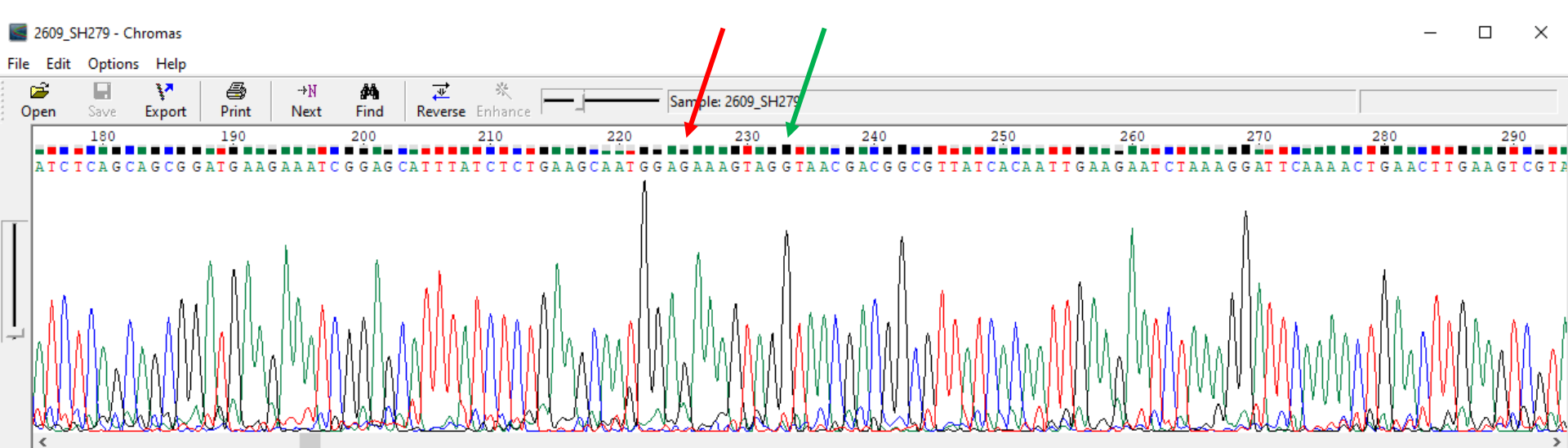
- Nanopore

General summary:

Mean read length:	386.9
Mean read quality:	11.8
Median read length:	99.0
Median read quality:	11.5
Number of reads:	11,266.0
Read length N50:	5,927.0
Total bases:	4,358,461.0

Hodnocení kvality stanovení báze

- Q score nebo Phred skóre
- Celočíselná hodnota představující odhadovanou pravděpodobnost chyby
- Kódování může být specifické dle platformy



Příklad hodnocení prostřednictvím Q-skóre u Sangerova sekvenování



PHRED skóre kvality

- Skóre kvality jsou reprezentována jako znaky ASCII
- Phred+33
 - v současnosti nejpoužívanější
 - Illumina 1.8+, Sanger, PacBio

! = nejnižší kvalita

~ = nejvyšší kvalita

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



- Phred+64
 - lze se s ním setkat u starších dat před r. 2018
 - Illumina 1.3+, Illumina 1.5+, Solexa
 - Před použitím dat může být nutný převod na Phred+33

Výpočet Phred Scores

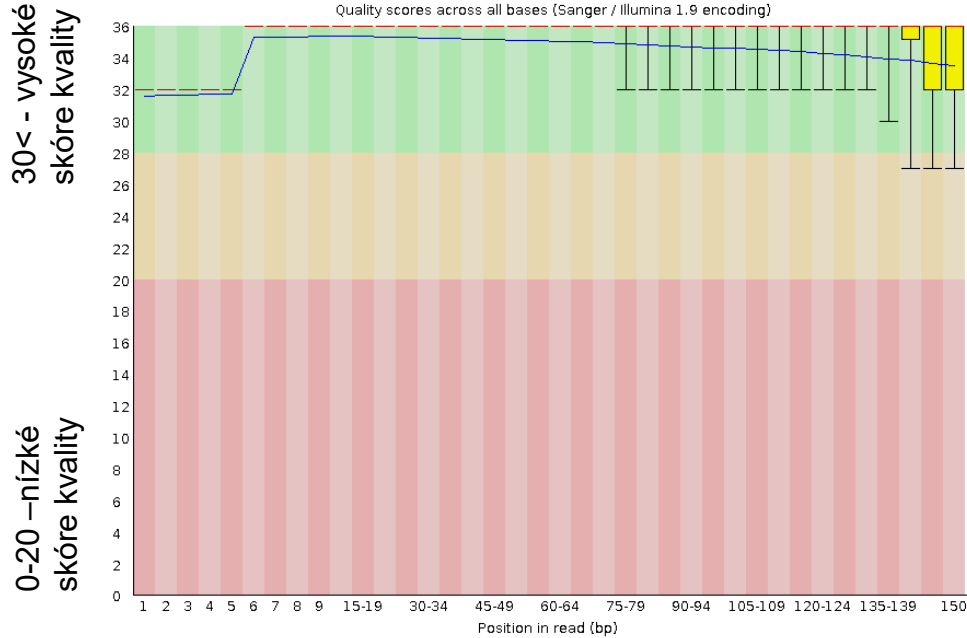
- Stanovení pravděpodobnosti, že uvedená báze je nesprávná: hodnota (p)
- Transformace na skóre Phred z plovoucí desetinné čárky, tedy kladné celé číslo:

$$\text{Phred} = -10 * (\text{int})\log_{10}(p)$$

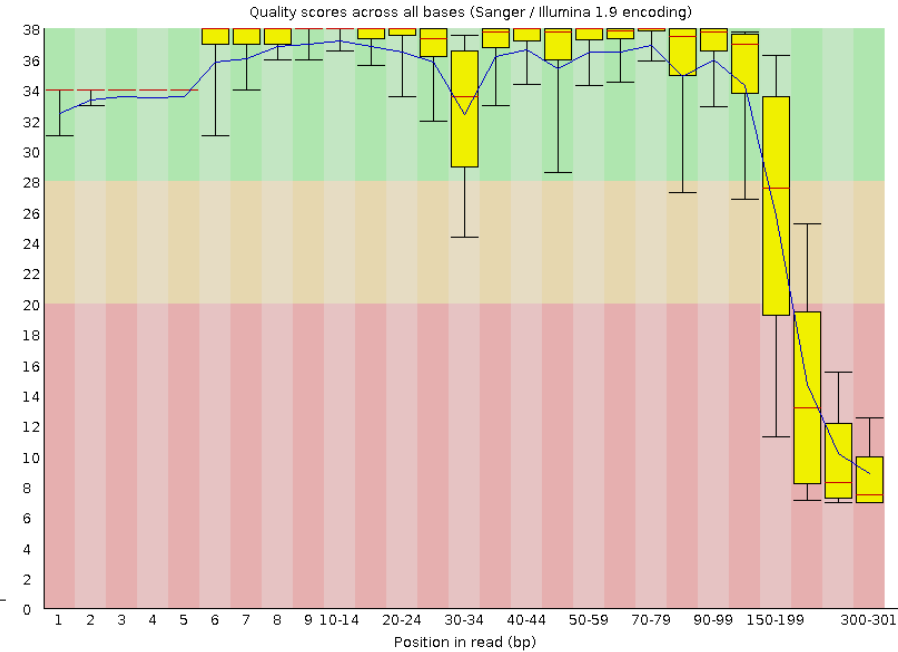
- p=0.1 Phred = 10
- p=0.01 Phred = 20
- p=0.001 Phred = 30

Statistické hodnocení kvality sekvence na bázi z Illumina čtení nástrojem FastQC

Dobrá kvalita

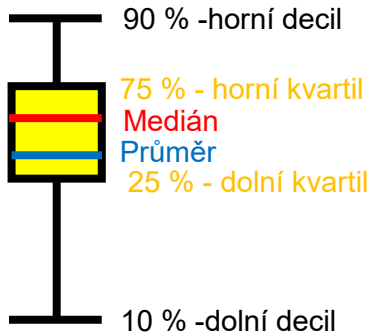


Špatná kvalita



metriky kontroly kvality, jsou hlášeny systémem varování na semaforech, normální (zelená), abnormální (oranžová), špatná (červená)

Distribuce skóre kvality

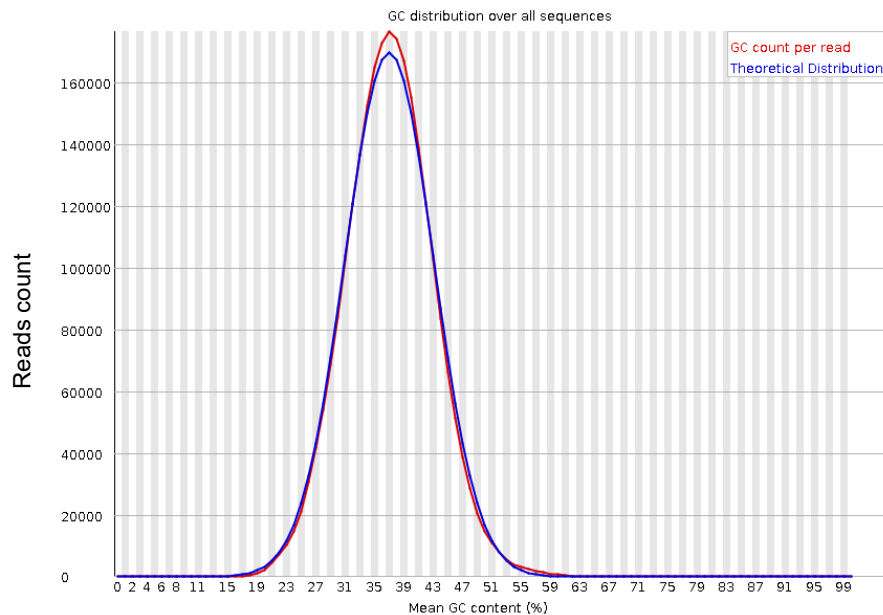


- [Basic Statistics](#)
- [Per base sequence quality](#)
- [Per tile sequence quality](#)
- [Per sequence quality scores](#)
- [Per base sequence content](#)
- [Per sequence GC content](#)
- [Per base N content](#)
- [Sequence Length Distribution](#)
- [Sequence Duplication Levels](#)
- [Overrepresented sequences](#)
- [Adapter Content](#)

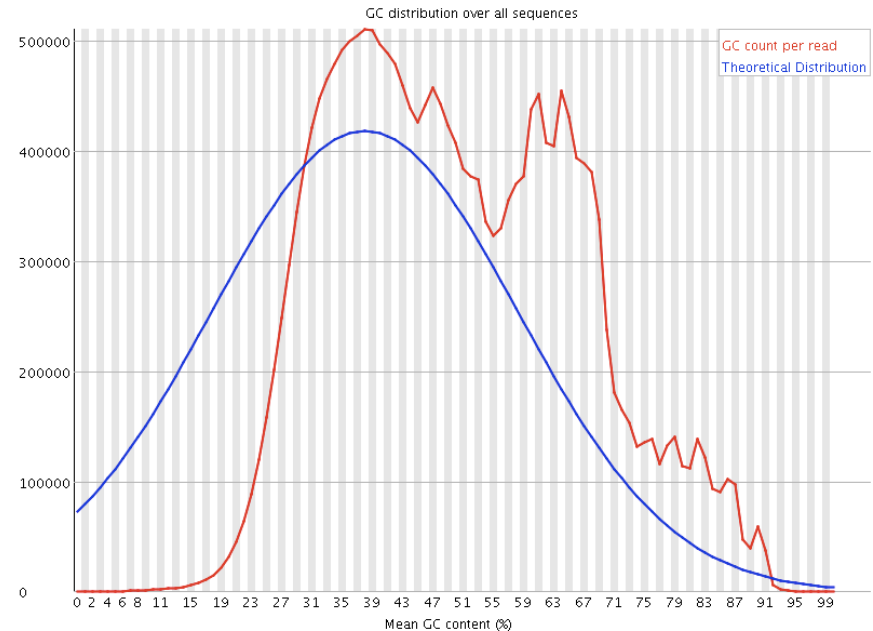
- [Basic Statistics](#)
- [Per base sequence quality](#)
- [Per sequence quality scores](#)
- [Per base sequence content](#)
- [Per base GC content](#)
- [Per sequence GC content](#)
- [Per base N content](#)
- [Sequence Length Distribution](#)
- [Sequence Duplication Levels](#)
- [Overrepresented sequences](#)
- [Kmer Content](#)

Hodnocení distribuce G+C pomocí FastQC

- Distribuce CG je první parametr, kde se zjistí kontaminace dat
- Zdroje kontaminace mohou být různé
 - technické sekvence DNA jako primery, adaptory, hostitelská DNA, rRNA
- Pokud je znám zdroj kontaminace, pak je možné tyto sekvence odfiltrovat



Distribuce obsahu CG pro každé čtení, je dobrá a pouze s jedním vrcholem. Data tedy pravděpodobně obsahují DNA pouze z jednoho organismu.



Nerovnoměrná distribuce CG % znamená pravděpodobnou kontaminaci. V tomto případě se jedná RNAseq s nedostatečně odstraněnou rRNA.

FastQC –kontaminace adaptory

⊗ Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCT	8122	8.122	Illumina Paired End PCR Primer 2 (100% over 40bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGATCGGAAG	5086	5.086	Illumina Paired End PCR Primer 2 (97% over 36bp)
AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTAC	1085	1.085	Illumina Single End PCR Primer 1 (100% over 40bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGGAAG	508	0.508	Illumina Paired End PCR Primer 2 (97% over 36bp)
AATTATACGGCGACCACCGAGATCTACACTCTTTCCCTAC	242	0.242	Illumina Single End PCR Primer 1 (97% over 40bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAAGATCGGAA	235	0.23500000000000001	Illumina Paired End Adapter 2 (96% over 31bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGATCGGAAGA	228	0.22799999999999998	Illumina Paired End Adapter 2 (96% over 28bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGGACG	205	0.20500000000000002	Illumina Paired End PCR Primer 2 (97% over 36bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGGATCGGAA	183	0.183	Illumina Paired End Adapter 2 (100% over 32bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGGTCGGAAG	183	0.183	Illumina Paired End Adapter 2 (100% over 32bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGAACT	164	0.164	Illumina Paired End PCR Primer 2 (97% over 40bp)
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGGTCT	129	0.129	Illumina Paired End PCR Primer 2 (97% over 40bp)
AATTATACTTCTACCACCTATATCTACACTCTTTCCCTAC	123	0.123	No Hit
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGGACT	122	0.122	Illumina Paired End PCR Primer 2 (97% over 36bp)
CGGTTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTTCAGC	113	0.11299999999999999	Illumina Paired End PCR Primer 2 (96% over 25bp)

Příliš mnoho adaptorových sekvencí může znamenat chybu při přípravě sekvenační knihovny. Řešením je opět odfiltrovat tyto sekvence, pokud je dostatek čtení. Pokud ne, může být nutné zopakovat sekvenování.

Filtrování podle kvality: co můžeme dělat s nekvalitními daty

- Odstranění nejhorších čtení podle skóre kvality
- Odstranění kontaminant
- Oříznutí (trimming)
 - Na začátku (adaptory, primery)
 - Na konci (klesající kvalita)
 - Automatické nebo manuální stanovení cutoff limitů
 - Při dostatečné hloubce pokrytí se upřednostňuje méně čtení o vyšší kvalitě
- Sloučení párových čtení
 - snížení redundance dat, pokud mají čtení velký překryv
- Obecně se podle typu dat a typu experimentu provádí odlišné kroky

Oříznutí čtení sekvence (trimming)

báze 1->5

zbytek adaptoru

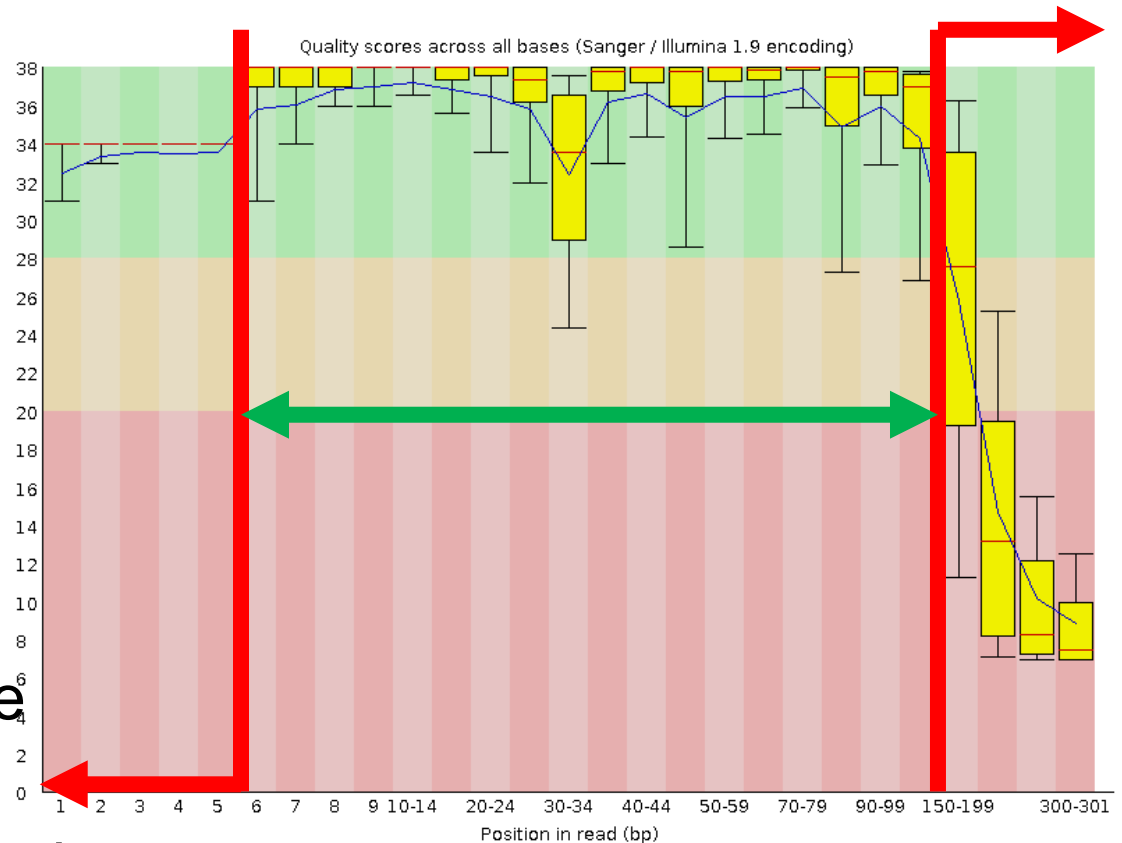
báze 6->150

zachovaný úsek

báze 150->300

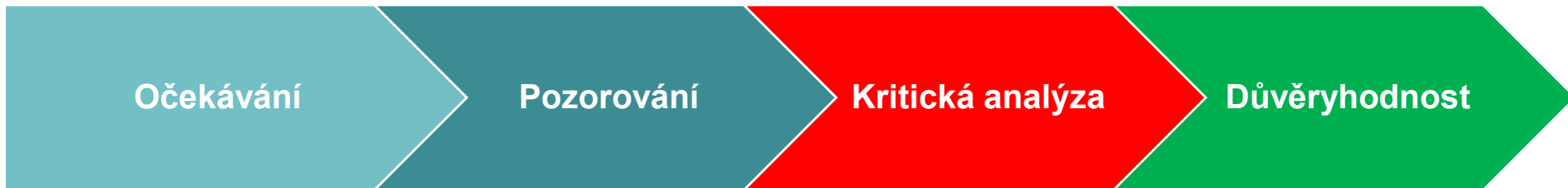
ztráta kvality

ořez cca od 200. báze



- cutadapt, trimmomatic
 - command line programy se širokými možnostmi filtrování a úprav čtení
- filolong
 - command line, pro filtrování kvality u dlouhých čtení

Závěr: Kontext je klíčem pro hodnocení kvality



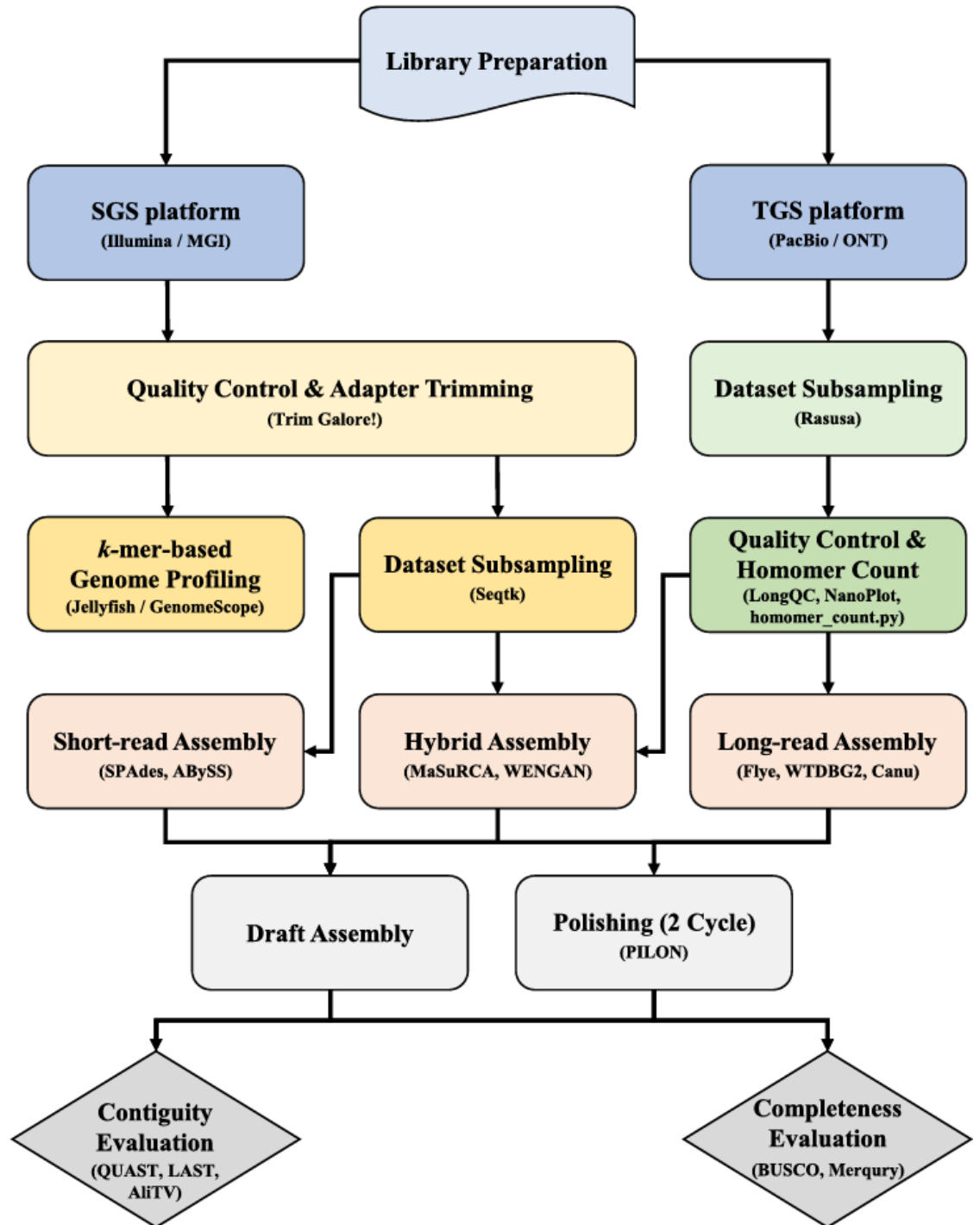
QC by měla být o tom, co očekáváte a co vidíte

Vývojový diagram sekvenování a assembly

Příklady programů

Závislé na:

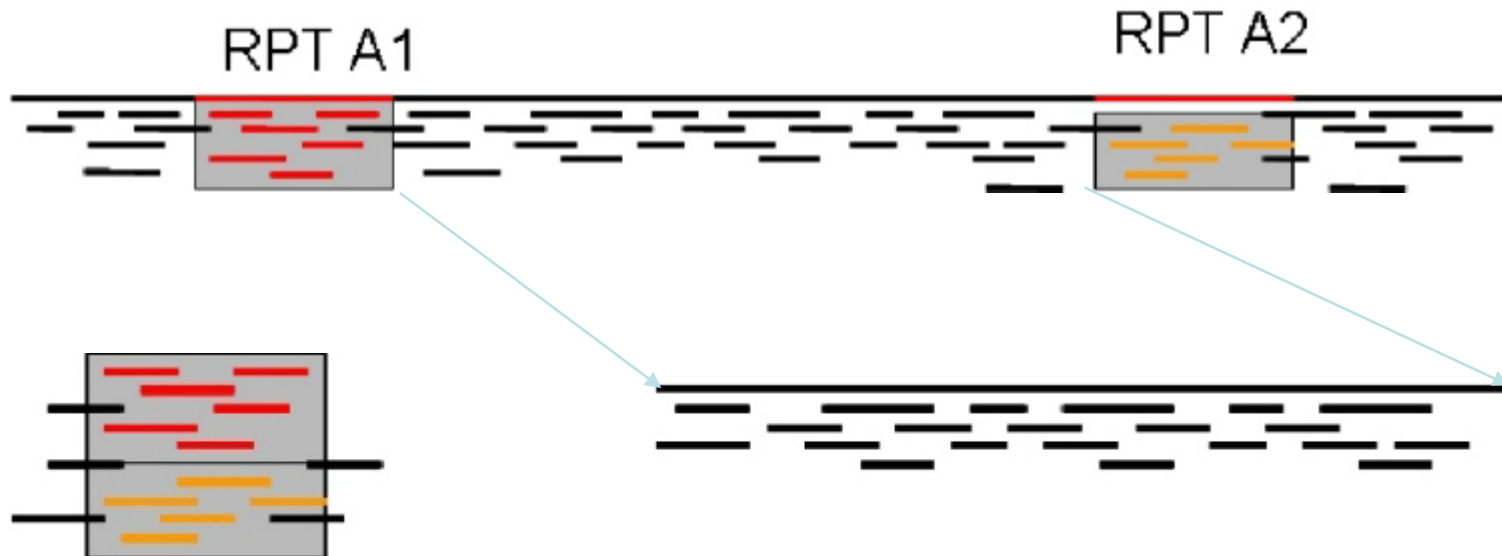
- Sekvenační platformě
- Získaném pokrytí
- Typu genomu



Repetice jsou příčinou rozdělení genomů do kontigů



Jestliže čtení je kratší než repetice → nemožnost sestavení sekvence



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

Kontig tvořený jedinečnou sekvencí, ohraničený repetitivními sekvencemi

- 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



Složení genomu – assembly

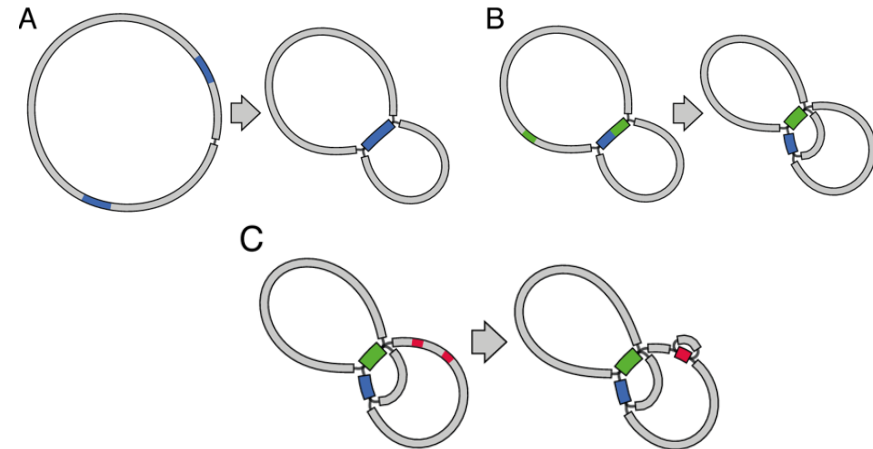
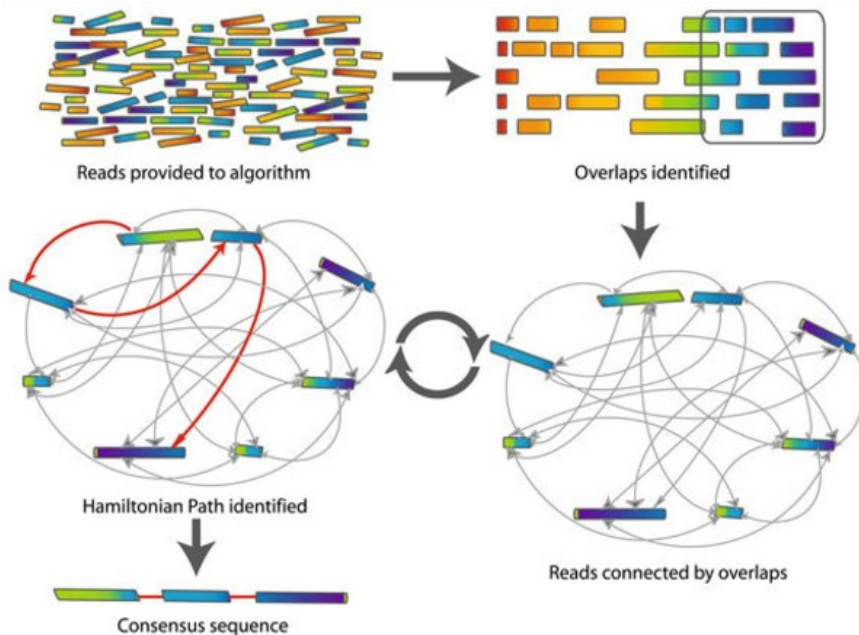
- Principy assembly

- OLC – overlap layout consensus

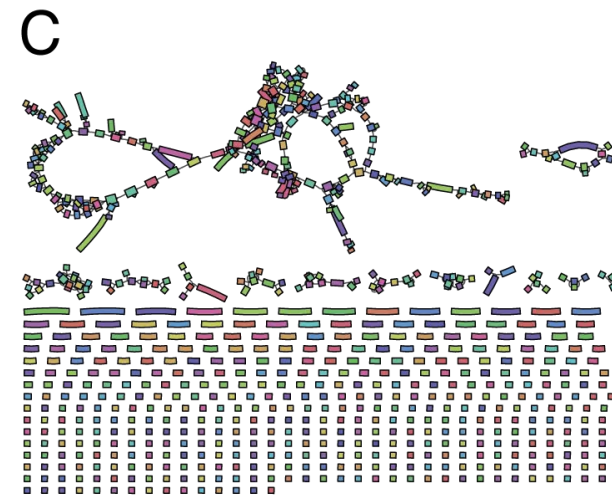
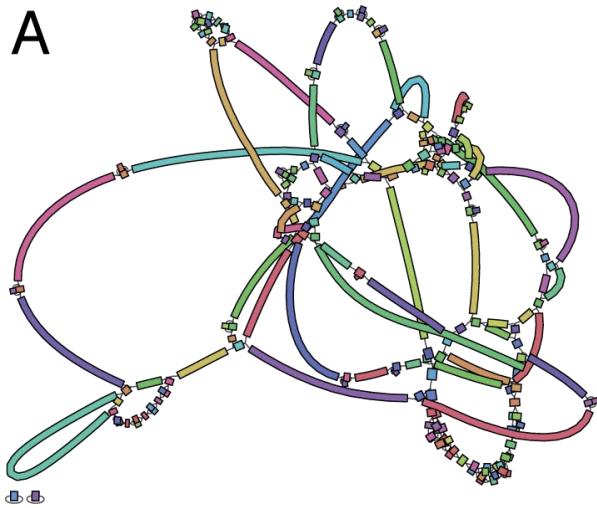
- výpočetně náročné
- vhodnější pro dlouhá čtení

- De Bruijnův graf

- ztráta části informace
- vhodnější pro krátká čtení



Příklad de Bruijnova grafu u mikrobiálního genomu (Illumina)



A – kvalitní sestavení

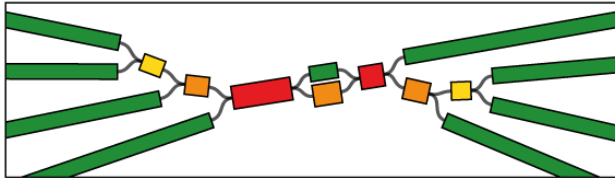
B – sestavení vyžadující optimalizaci, kombinace dlouhých a krátkých kontigů

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

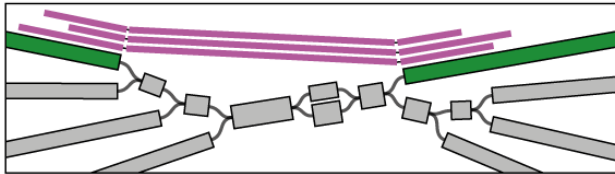
Hybridní assembly a bridging

- Kombinace krátkých čtení (Illumina, IonTorrent) a dlouhých čtení (PacBio, Nanopore) umožňuje hybridní assembly
 - repetice vyřešeny namapováním dlouhých čtení
 - chyby vyřešeny namapováním krátkých čtení

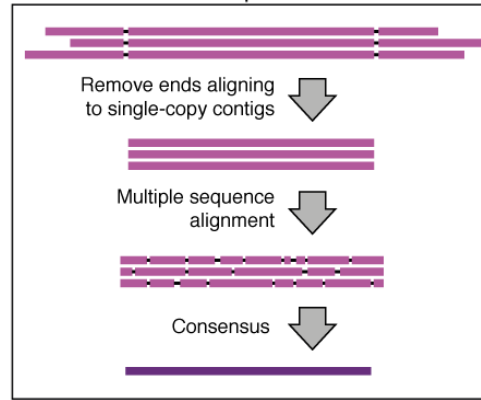
Repeat region in unbridged graph



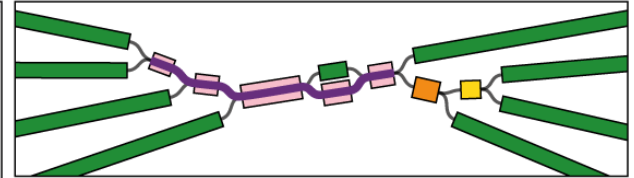
Semi-global long read alignment



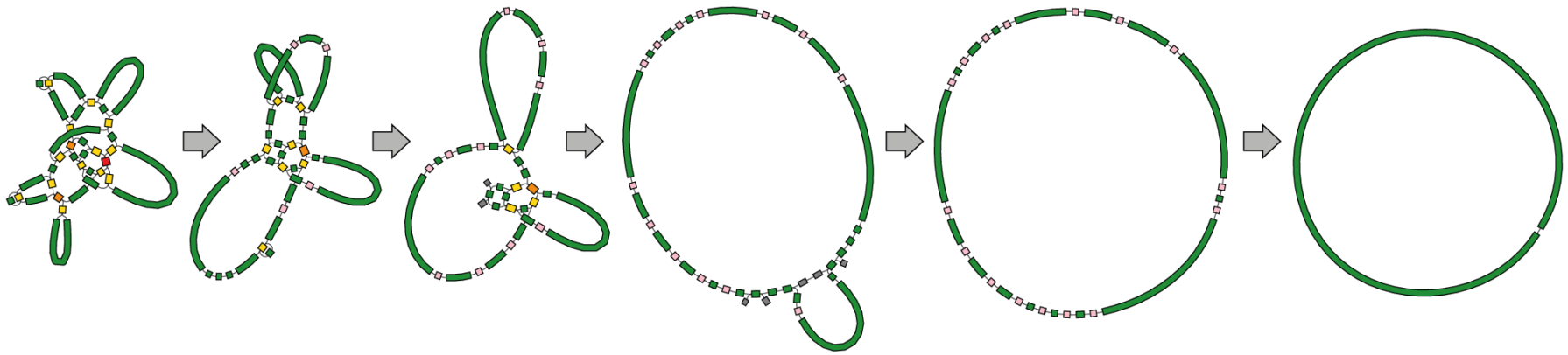
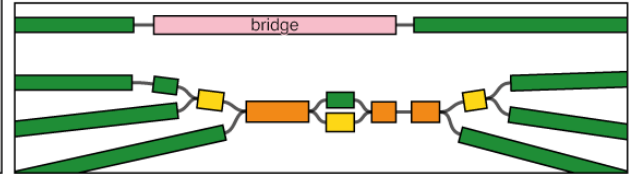
Consensus read sequence



Path finding



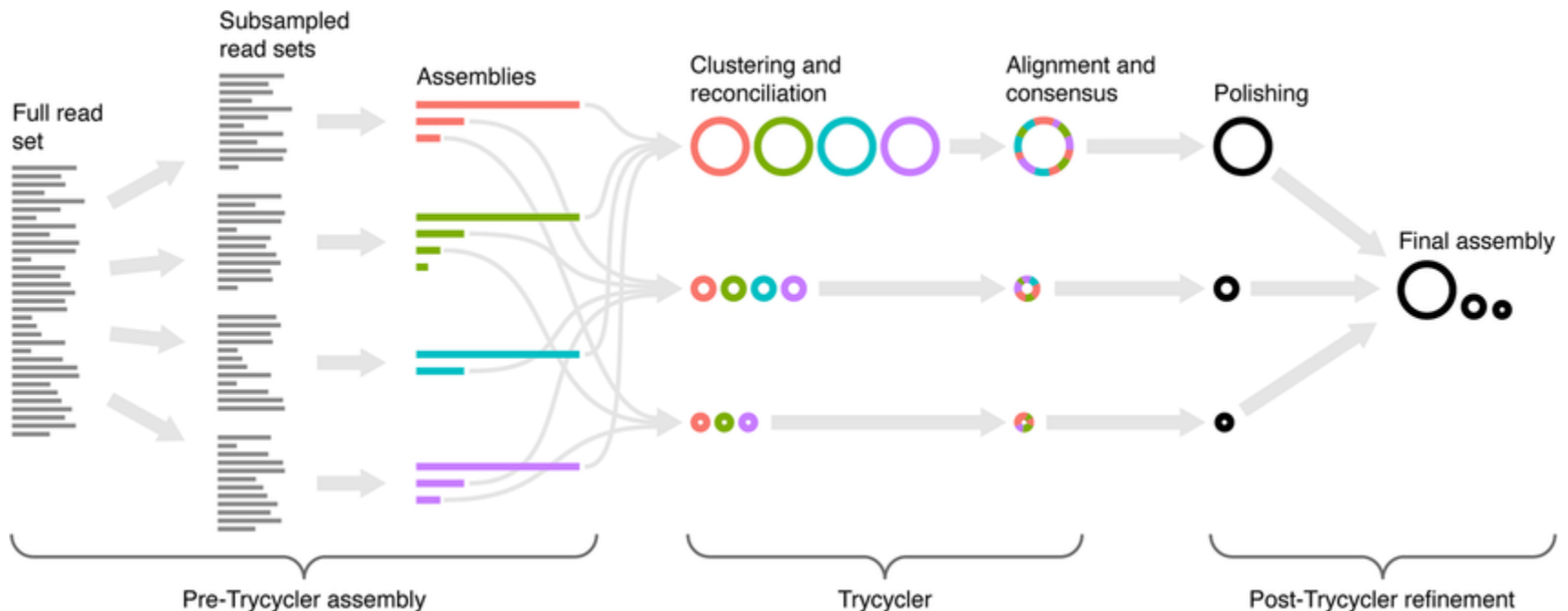
Bridged graph



Pro assembly obvykle využíváme připravené pipeline

Mohou být založeny na kombinaci různých assemblerů

- Unicycler: Short reads and hybrid assembly pipeline
- Tricycler: long-read assembly pipeline



Klasifikace assemblerů a využívané strategie

Assembler	Error correction	Contig extension	Complexity reduction
Flye	Direct	Graph-based disjointing correction	Disjointig construction
WTDBG2		Fuzzy Bruijn graph	Hash table based on the k -mer block (bin)
		Partial order alignment	Non-redundant k -mer removing
Canu	Hierarchical	Best overlap graph	tf -idf weighted MHAP
		CABOG	
		PBDAG-CON	
MaSuRCA	Hybrid	de Bruijn	Super-read construction
		CABOG	
WENGAN		de Bruijn	
		Partial order alignment	Synthetic scaffolding graph
SPAdes	Short-read only	Multisized de Bruijn	—
		de Bruijn	—
ABYSS		Paired-end-based contig extension	—



Statistiky assembly

- N50 a L50 statistiky sady délek kontigů nebo skafoldů, prostřednictvím kterých můžeme srovnat kvalitu assembly
- **N50** definuje kvalitu assembly z hlediska spojitosti. **Délka** sekvence nejkratšího kontigu který přispěl k sestavení 50 % celkové délky.
- **L50** je definován jako nejmenší **počet** kontigů, jejichž součet délek tvoří polovinu velikosti genomu.



Automatická anotace

- Algoritmy- viz lekce Hledání genů
- Servery pro automatickou anotaci
 - **NCBI Prokaryotic Genome Annotation Pipeline (PGAP)**
https://www.ncbi.nlm.nih.gov/genome/annotation_prok/process/
 - **NCBI Eukaryotic Genome Annotation Pipeline**
https://www.ncbi.nlm.nih.gov/genome/annotation_euk/process/
 - **RAST**
<http://rast.nmpdr.org/>
- Offline command line aplikace
 - **Prokka**, bakteriální genomy
<https://github.com/tseemann/prokka>
 - **MARK**, eukaryotické genomy
<https://reslp.github.io/blog/My-MAKER-Pipeline/>

Genome

Genome ▾

[Limits](#) [Advanced](#)

Search

Prokaryotic Annotation Home

Documentation ▾

Complete Genome Submission ▾

WGS Genome Submission ▾

NCBI Prokaryotic Genome Annotation Pipeline

NCBI Prokaryotic Genome Annotation Pipeline (PGAP) is designed to annotate bacterial and archaeal genomes (chromosomes and plasmids).

Genome annotation is a multi-level process that includes prediction of protein-coding genes, as well as other functional genome units such as structural RNAs, tRNAs, small RNAs, pseudogenes, control regions, direct and inverted repeats, insertion sequences, transposons and other mobile elements.

NCBI has developed an automatic prokaryotic genome annotation pipeline that combines *ab initio* gene prediction algorithms with homology based methods. The first version of NCBI Prokaryotic Genome Pipeline was developed in 2001 and is regularly upgraded to improve structural and functional annotation quality ([Haft DH et al 2018](#), [Tatusova T et al 2016](#)). Recent improvements utilize curated protein profile hidden Markov models (HMMs), including [TIGRFAMS](#) and new HMMs for antimicrobial resistance proteins, and curated complex domain architectures for functional annotation of proteins. NCBI's annotation pipeline depends on several internal databases and is not currently available for download or use outside of the NCBI environment.

Related documentation:

- [Annotation process](#)
- [Annotation standards](#)
- [Assemblies excluded from RefSeq](#)
- [Release notes](#)

GenBank

The NCBI prokaryotic annotation pipeline is available as a service for GenBank submitters. The pipeline is capable of annotating both complete genomes and draft WGS genomes consisting of multiple contigs. You can request PGAP annotation when you submit your genome to GenBank.

Both WGS and non-WGS genomes, including gapless complete bacterial chromosomes, can be submitted via the Submission Portal. You will be asked to choose whether the genome being submitted is considered WGS or not. The differences for GenBank purposes are:

non-WGS:

- Each chromosome is in a single sequence and there are no extra sequences
- Each sequence in the genome must be assigned to a chromosome or plasmid or organelle
- Plasmids and organelles can still be in multiple pieces.

WGS:

- One or more chromosomes are in multiple pieces and/or some sequences are not assembled into chromosomes

Genome

Genome ▾

[Limits](#) [Advanced](#)

Search

[Eukaryotic Annotation Home](#)[Documentation](#) ▾[Annotated Genomes](#) ▾[Annotation Policy](#)[Request Annotation](#)

The NCBI Eukaryotic Genome Annotation Pipeline

The NCBI Eukaryotic Genome Annotation Pipeline provides content for various NCBI resources including [Nucleotide](#), [Protein](#), [BLAST](#), [Gene](#) and the [Genome Data Viewer](#) genome browser.

This page provides an overview of the annotation process. Please refer to [the Eukaryotic Genome Annotation chapter of the NCBI Handbook](#) for algorithmic details.

The pipeline uses a modular framework for the execution of all annotation tasks from the fetching of raw and curated data from public repositories (sequence and [Assembly](#) databases) to the alignment of sequences and the prediction of genes, to the submission of the accessioned annotation products to public databases. Core components of the pipeline are alignment programs ([Splign](#) and [ProSplign](#)) and an HMM-based gene prediction program ([Gnomon](#)) developed at NCBI.

Important features of the pipeline include:

- flexibility and speed
- higher weight given to curated evidence than non-curated evidence
- utilization of RNA-Seq for gene prediction
- production of models that compensate for assembly issues
- tracking of gene loci from one annotation to the next
- ability to co-annotate multiple assemblies for the same organism

The products of an annotation run (chromosome, scaffolds and model transcripts and proteins) are labeled with an Annotation Release number. The Annotation Release name is the combination of the organism name and Annotation Release number (e.g. NCBI *Pongo abelii* Annotation Release 103) and is used throughout NCBI as a way to uniquely identify annotation products originating from the same annotation run.

Contents

- [Process](#)
 - [Source of genome assemblies](#)
 - [Masking](#)
 - [Transcript alignments](#)
 - [RNA-Seq read alignments](#)
 - [Protein alignments](#)
 - [Model prediction](#)
 - [Curated RefSeq genomic sequence alignments](#)
 - [Choosing the best models for a gene](#)
 - [Protein naming and determination of locus type](#)
 - [Assignment of GeneIDs](#)
 - [Annotation of small RNAs](#)

RAST (Rapid Annotation using Subsystem Technology) Server

<http://rast.nmpdr.org/>

Upload a Genome

Complete Upload

Please consider the following options for the RAST annotation pipeline:

RAST Annotation Settings:

Choose RAST annotation scheme: Choose "Classic RAST" for the current production RAST, or "RASTtk" for the new modular RAST pipeline currently in testing.

Customize RASTtk pipeline: Yes Customize the RASTtk pipeline

Stage name	Enabled	Parameters	Condition
‡ call-features-rRNA-SEED	<input checked="" type="checkbox"/> Yes		
‡ call-features-tRNA-trnscan	<input checked="" type="checkbox"/> Yes		
‡ call-features-repeat-region-SEED	<input checked="" type="checkbox"/> Yes	Minimum identity 95 Minimum length 100	
‡ call-selenoproteins	<input checked="" type="checkbox"/> Yes		
‡ call-pyrrolysoproteins	<input checked="" type="checkbox"/> Yes		
‡ call-features-insertion-sequences	<input type="checkbox"/> Yes		
‡ call-features-strep-suis-repeat	<input checked="" type="checkbox"/> Yes		\$genome->{scientific_name}
‡ call-features-strep-pneumo-repeat	<input checked="" type="checkbox"/> Yes		\$genome->{scientific_name}
‡ call-features-crispr	<input checked="" type="checkbox"/> Yes		
‡ call-features-CDS-glimmer3	<input checked="" type="checkbox"/> Yes	Minimum training length 2000	
‡ call-features-CDS-prodigal	<input checked="" type="checkbox"/> Yes		
‡ call-features-CDS-genemark	<input type="checkbox"/> Yes		
‡ annotate-proteins-kmer-v2	<input checked="" type="checkbox"/> Yes	Minimum kmer hits required 5 Only annotate hypothetical proteins <input type="checkbox"/> Yes	
‡ annotate-proteins-kmer-v1	<input checked="" type="checkbox"/> Yes	Kmer dataset to use Release70 <input type="radio"/> Release59 <input type="radio"/>	
‡ annotate-proteins-similarity	<input checked="" type="checkbox"/> Yes	Only annotate hypothetical proteins <input checked="" type="checkbox"/> Yes	
‡ resolve-overlapping-features	<input checked="" type="checkbox"/> Yes		
‡ find-close-neighbors	<input checked="" type="checkbox"/> Yes		
‡ call-features-prophage-phispy	<input type="checkbox"/> Yes		

Automatically fix errors? Yes *The automatic annotation process may run into problems, such as gene candidates overlapping RNAs, or genes embedded inside other genes. To automatically resolve these problems (even if that requires deleting some gene candidates), please check this box.*

Fix frameshifts? Yes *If you wish for the pipeline to fix frameshifts, check this option. Otherwise frameshifts will not be corrected.*

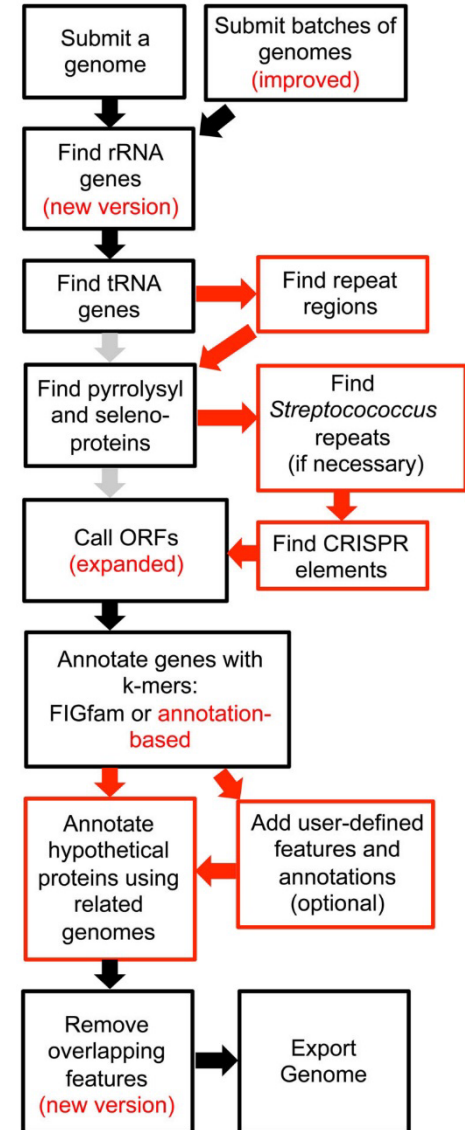
Build metabolic model? Yes *If you wish RAST to build a metabolic model for this genome, check this option.*

Turn on debug? Yes *If you wish debug statements to be printed for this job, check this box.*

Set verbose level 0 *Set this to the verbosity level of choice for error messages.*

Disable replication Yes *Even if this job is identical to a previous job, run it from scratch.*

Finish the upload

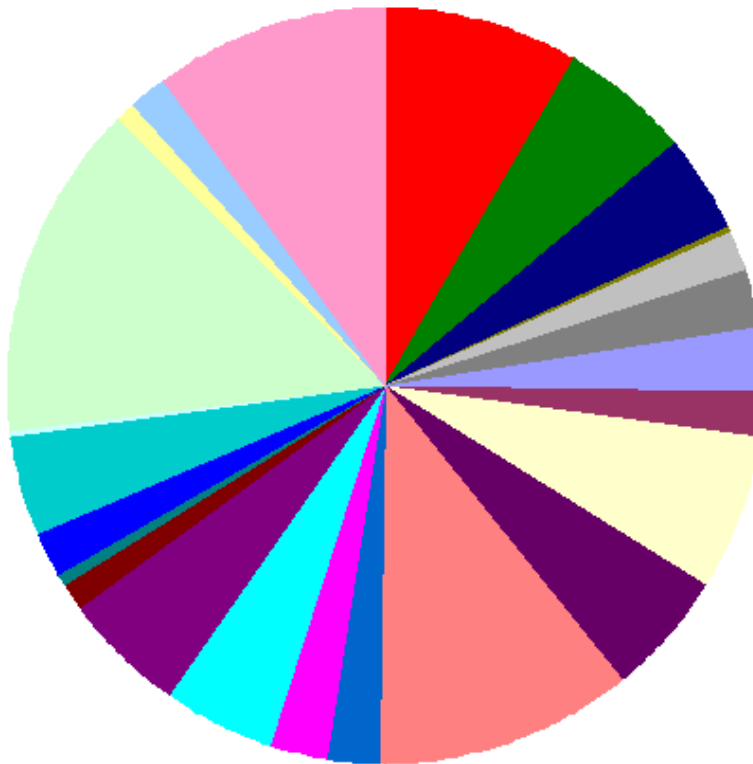


RAST anotace bakteriálního genomu a klasifikace do subsystemů

Subsystem Coverage



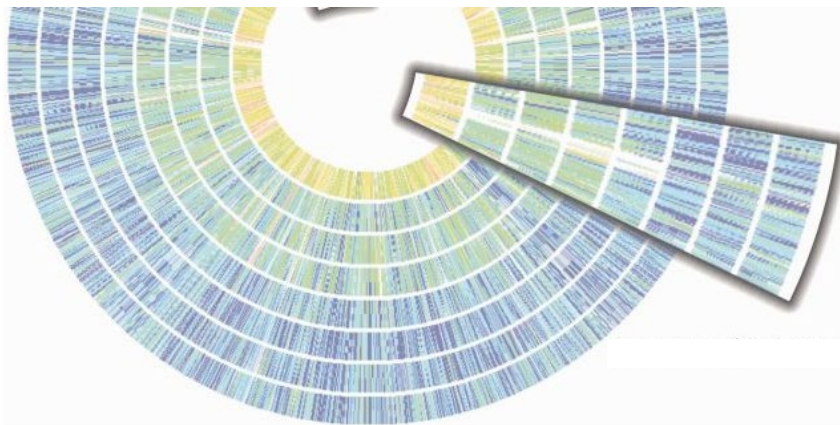
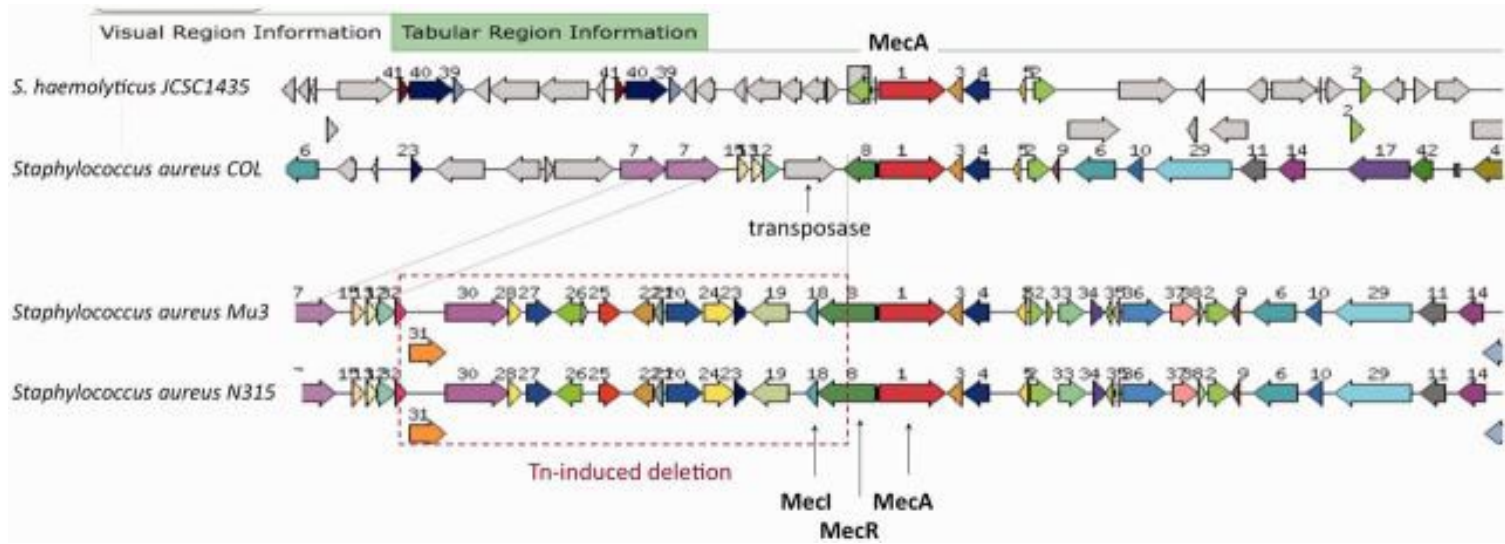
Subsystem Category Distribution



Subsystem Feature Counts

- ⊕ Cofactors, Vitamins, Prosthetic Groups, Pigments (152)
- ⊕ Cell Wall and Capsule (100)
- ⊕ Virulence, Disease and Defense (71)
- ⊕ Potassium metabolism (6)
- ⊕ Photosynthesis (0)
- ⊕ Miscellaneous (32)
- ⊕ Phages, Prophages, Transposable elements, Plasmids (45)
- ⊕ Membrane Transport (50)
- ⊕ Iron acquisition and metabolism (32)
- ⊕ RNA Metabolism (122)
- ⊕ Nucleosides and Nucleotides (90)
- ⊕ Protein Metabolism (200)
- ⊕ Cell Division and Cell Cycle (40)
- ⊕ Motility and Chemotaxis (0)
- ⊕ Regulation and Cell signaling (46)
- ⊕ Secondary Metabolism (4)
- ⊕ DNA Metabolism (82)
- ⊕ Fatty Acids, Lipids, and Isoprenoids (93)
- ⊕ Nitrogen Metabolism (23)
- ⊕ Dormancy and Sporulation (11)
- ⊕ Respiration (31)
- ⊕ Stress Response (74)
- ⊕ Metabolism of Aromatic Compounds (5)
- ⊕ Amino Acids and Derivatives (262)
- ⊕ Sulfur Metabolism (15)
- ⊕ Phosphorus Metabolism (27)
- ⊕ Carbohydrates (173)

Vizualizace úseku anotovaného genomu a srovnávací analýzy



Vizualizace lokusů:

- Srovnání anotace genů

Vizualizace celých genomů:

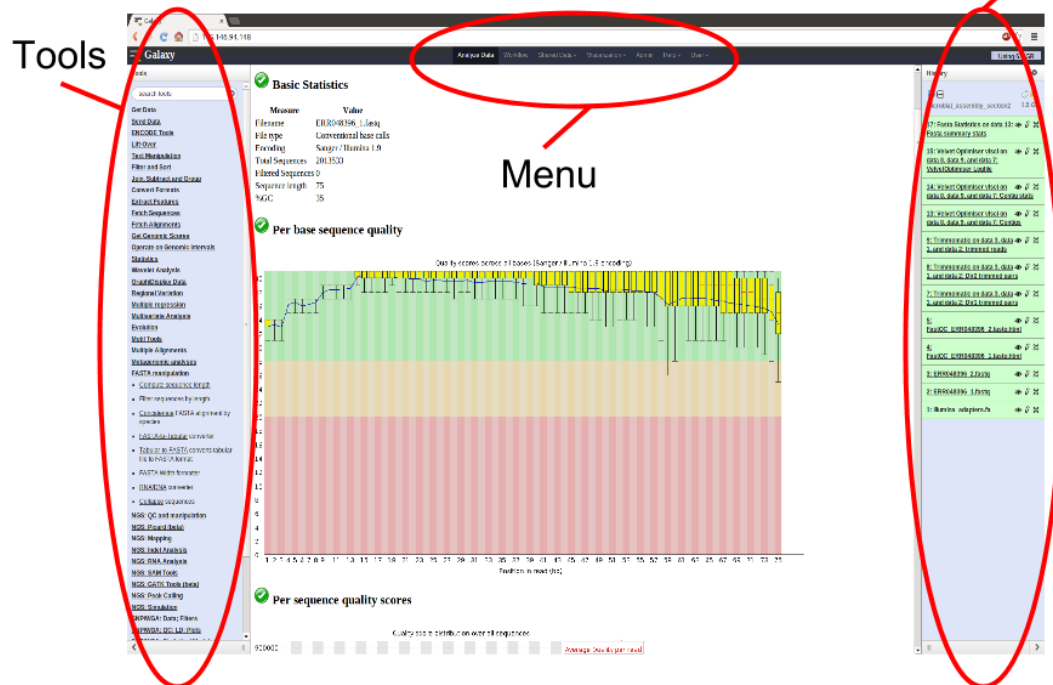
- Míra konzervovanosti
- Inzerce/ delece

Platformy pro manipulaci sekvencí



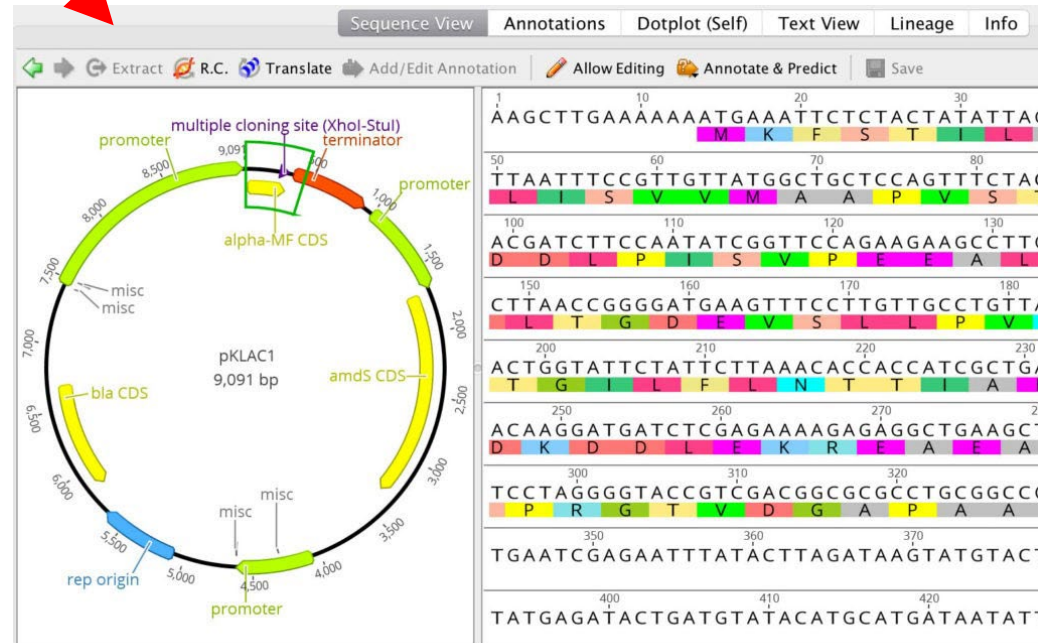
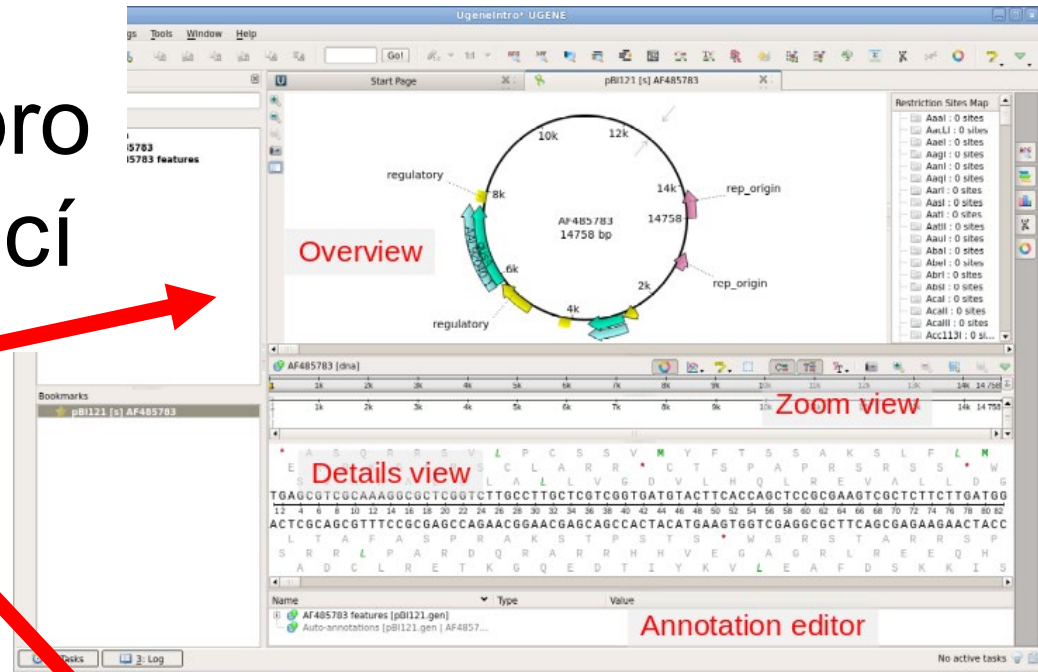
- usnadňují práci uživatelům, kteří nejsou zblhlí v programování a administraci systémů
- Webový server **Galaxy**
 - open source
 - připravené workflows
 - od jednoduchých textových manipulací po analýzu sekvenčních dat, genome assembly, metagenomiku
 - nástroje pro vizualizaci dat
 - <https://usegalaxy.org>

Analysis History

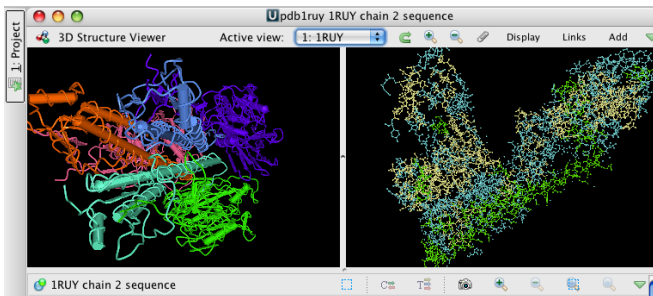


Lokální platformy pro manipulaci sekvencí

- neplacené (např. **Ugene**) nebo licencované (**Geneious**)
- příjemné grafické rozhraní
- implementace nástrojů
 - anotace DNA, proteinů
 - prohledávání databází
 - multiple sequence alignment
 - PCR – návrh primerů
 - contig assembly
 - short reads alignment
 - search ORF
 - 3D structure Viewer
 - ...



Alt click on a sequence position or annotation, or select a region to zoom in. Alt-shift click to zoom out.



Jak se data dostanou do databází?

Bacteria and Bacteriophage

A Compilation from the Genbank® and EMBL Data Libraries

Book 1087

PHIX174CS: bacteriophage phix174 gene a protein cleavage site. [ss-DNA]
EMBL ID: PH174CS ACCESSION NUMBERS: K00813 DATE: pre-entry 84-06-01
REFERENCES: [1] (bases 1 to 65) Brown,D.R., Schmidt-Glenewinkel,T., Reinberg,D. and Hurwitz,J.; "dna sequences which support activities of the bacteriophage phix174 gene a protein"; J Biol Chem 258, 8402-8412 (1983)
KEYWORDS: closing activity; gene A; nicking activity; origin of replication; replication form.
SOURCE: bacteriophage phix174 (strain am3) dna, grown in e.coli hf4704. Bacteriophage phi-X174
COMMENT: the gene a protein of phix174 mediates both initiation and termination of viral strand dna synthesis. it nicks and closes within a 30 nucleotide sequence (bases 13-42) which is well conserved among icosohedral ss(c) dna phages; this 30 bp region is implicated as the specific target for gene a protein action, and as the essential dna sequence required for replication origin function.
SITES: key site span description
refnumbr 1 1 sequence not numbered in [1]
cutss 20 0 gene a protein cleavage site
ORIGIN: replication origin of phix174.
SEQUENCE: 65 bp 23 a 18 c 12 g 12 t
1 aatgtgctcc cccaacttga tattaataac actatagacc accgccccga aggggacgaa aatg

PHIX174DE: bacteriophage phix174 d and e genes.
EMBL ID: * ACCESSION NUMBERS: J02483 DATE: pre-entry 83-03-01
REFERENCES: [1] (bases 1 to 521) Barrell,B.G., Air,G.M. and Hutchison,C.A.III.; "overlapping genes in bacteriophage phix174"; Nature 264, 34-41 (1976)
SOURCE: phix174. Bacteriophage phi-X174
SEQUENCE: 521 bp 125 a 115 c 124 g 157 t
1 gagtcgatg ctgtcaacc actaataggt aagaaatcat gagtcaagtt actgaacaat cgtacgttt ccagacogct ttggcctcta ttaagctcat
101 tcaggcttct gcggttttgg atttaaccga agatgatitc gattitctga cgagtaacaa agtttgatt gctactgacc gctctcgtgc tctcgtcgc
201 gttgaggctt gcgtttatgg tacgctggac ttgtgggat accctcgott tctgctcct gttgagitta ttgctgocgt cattgcttat tatgttcato
301 cgtcaacat tcaaacggcc tctctcatca tggaaaggocg tgaatttacg gaaaacatta ttaatggcgt cgagcgtccg gttaaagccg ctgaattgtt
401 cgcgtttacc ttgcgtgtac ggcaggaata cactgacgtt ctactgacg cagaagaaaa cgtgcgtcaa aaattacgtg cggaaggagt gatgtaatg
501 ctaaaggtaa aaaacgttct g

Jak se data dostanou do databází?

- **Submission - Vložení do databáze**
- **webový portál**
 - **BankIt** (GenBank)
<https://submit.ncbi.nlm.nih.gov/about/bankit/>
 - **Submission Portal** (GenBank)
<https://submit.ncbi.nlm.nih.gov/subs/genome/>
 - **WebIn** (EMBL/European Nucleotide Archive)
<http://www.ebi.ac.uk/ena/submit>
 - **Sakura** (DDBJ)
<http://www.ddbj.nig.ac.jp/sub/websub-e.html>
- **samostatná aplikace pro PC**
 - **Sequin**, delší manuálně anotované sekvence, fylogenetické, populační nebo mutační studie obsahující sekvenční přílohy
http://www.ncbi.nlm.nih.gov/Sequin/download/seq_download.html
 - **Tbl2asn**, command line program, celé genomy EST, STS a zasílání velkých dávek sekvencí, automatizuje vytvoření záznamu sekvence
- **Minimální požadavky pro vložení**

GenBank formát

- Záznam anotovaného genomu
- Skládá se ze tří sekcí
 - **Header**
Informace o vlastnostech sekvence a jejím zdroji
 - **Feature Table**
Anotacemi formou deskriptorů; u genů může obsahovat i jejich translaci do proteinu
 - **Sequence**
Vlastní nukleotidová sekvence
- Historické omezení na 60 znaků na řádek
 - pole v hlavičce tak mají maximální délku
 - v současnosti už neplatí striktně

Tradiční záznam GenBank

```
LOCUS      AY182241                1931 bp    mRNA    linear    PLN 04-MAY-2004
DEFINITION Malus x domestica (E,E)-alpha-farnesene synthase (AFS1) mRNA,
            complete cds.
ACCESSION  AY182241
VERSION    AY182241.2  GI:32265057
KEYWORDS   .
SOURCE     Malus x domestica (cultivated apple)
ORGANISM   Malus x domestica
            Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
            Spermatophyta; Magnoliophyta; eudicotyledons; core eudicots;
            rosids; eurosids I; Rosales; Rosaceae; Maloideae; Malus.
REFERENCE  1 (bases 1 to 1931)
AUTHORS    Pechous,S.W. and Whitaker,B.D.
TITLE      Cloning and functional expression of an (E,E)-alpha-farnesene
            synthase cDNA from peel tissue of apple fruit
JOURNAL    Planta 219, 84-94 (2004)
REFERENCE  2 (bases 1 to 1931)
AUTHORS    Pechous,S.W. and Whitaker,B.D.
TITLE      Direct Submission
JOURNAL    Submitted (18-NOV-2002) PSI-Produce Quality and Safety Lab,
            USDA-ARS, 10300 Baltimore Ave. Bldg. 002, Rm. 205, Beltsville, MD
            20705, USA
REFERENCE  3 (bases 1 to 1931)
AUTHORS    Pechous,S.W. and Whitaker,B.D.
TITLE      Direct Submission
JOURNAL    Submitted (25-JUN-2003) PSI-Produce Quality and Safety Lab,
            USDA-ARS, 10300 Baltimore Ave. Bldg. 002, Rm. 205, Beltsville, MD
            20705, USA
REMARK     Sequence update by submitter
COMMENT    On Jun 26, 2003 this sequence version replaced gi:27804758.
FEATURES   Location/Qualifiers
            source          1..1931
                        /organism="Malus x domestica"
                        /mol_type="mRNA"
                        /cultivar="'Law Rome'"
                        /db_xref="taxon:3750"
                        /tissue_type="peel"
            gene            1..1931
                        /gene="AFS1"
            CDS             54..1784
                        /gene="AFS1"
                        /note="terpene synthase"
                        /codon_start=1
                        /product="(E,E)-alpha-farnesene synthase"
                        /protein_id="AAO22848.2"
                        /db_xref="GI:32265058"
                        /translation="MEFRVHLQADNEQKIFQNQMKPEPEASYLINQRRSANYKPNWIK
            NDFLDQSLISKYDGDYRKLSEKLIIEVKIYISAETMDLVAKLELIDSVRKGLANLF
            EKEIKALDSIAAIESDNLGTRDDLYGTALHFKILRQHGYKVSQDIFGRFMDEKGTLE
            DFLHKNEDLLYINSLIVRLNNDLGTSAAEQERGDSPSSIVCYMREVNASEETARKNIK
            GMIDNAWKVNGKCFITTQVPLSSFMNNATNMARVAHSLYKDGDFGQEKGRPTHI
            LSLLFQPLVN"
ORIGIN
1  ttctgtatc  ccaaacatct  cgagcttctt  gtacaccaa  ttaggtattc  actatggaat
61  tcagagttca  cttgcaagct  gataatgagc  agaaaatttt  tcaaaaccag  atgaaacccg
121  aacctgaagc  ctcttacttg  attaatacaa  gacggctctg  aaattacaag  ccaaatattt
181  ggaagaacga  tttcctagat  caatctctta  tcagcaaata  cgatggagat  gagtatogga
241  agctgtctga  gaagttaata  gaagaagtta  agatttatat  atctgctgaa  acaatggatt
//
```

Header

Feature Table

Sequence

Identifikace záznamu v primárních sekvenčních databázích

• Databáze

- GenBank
- EMBL-Bank (European Nucleotide Archive, ENA)
- DDBJ

• Identifikátory sekvence

GenBank ▾

**Escherichia coli citrate lyase beta-subunit (citE) gene, partial cds
lyase gamma-subunit (citD), citrate lyase ligase (citC), histidine kinase
two component response regulator (dpiA) genes, complete cds**

GenBank: U46667.1

[FASTA](#) [Graphics](#)

Go to:

LOCUS	ECU46667	5024 bp	DNA	linear	BCT 25-JUL-2016
DEFINITION	Escherichia coli citrate lyase beta-subunit (citE) gene, partial cds; and citrate lyase gamma-subunit (citD), citrate lyase ligase (citC), histidine kinase (dpiB), and two component response regulator (dpiA) genes, complete cds.				
ACCESSION	U46667				
VERSION	U46667.1				

Send to: ▾

Complete Record
 Gene Features

Choose Destination

File Clipboard
 Collections Analysis Tool

Download 1 item.

Format
GenBank ▾

Show GI

Create File

Related information
Protein



Historie verzí

- Sequence Revision History tool
 - Struktura zápisu:

<http://www.ncbi.nlm.nih.gov/nuccore/U46667?report=girevhist>

Revision History ▾

Send to: ▾

Show difference between **I** and **II** as

[Escherichia coli citrate lyase beta-subunit \(citE\), gene, partial cds; and citrate lyase gamma-subunit \(citD\), citrate lyase ligase \(citC\), histidine kinase \(dpiB\), and two component response regulator \(dpiA\) genes, complete cds](#)

5,024 bp linear DNA

Accession: U46667.1 GI: 3172140

Current status: live

I	II	Version	Gi	Accession	Update Date	Action
<input checked="" type="radio"/>	<input type="radio"/>	1	3172140	U46667.1	Jul 25, 2016 12:54 PM	
<input type="radio"/>	<input checked="" type="radio"/>	1	3172140	U46667.1	Jun 23, 2010 09:27 AM	
<input type="radio"/>	<input type="radio"/>	1	3172140	U46667.1	Nov 30, 2009 01:55 PM	
<input type="radio"/>	<input type="radio"/>	1	3172140	U46667.1	Aug 7, 1998 09:28 AM	
<input type="radio"/>	<input type="radio"/>	1	3172140	U46667.1	Jun 2, 1998 04:31 PM	
<input type="radio"/>	<input type="radio"/>	0	2734632	U46667.0	Jan 3, 1998 12:12 AM	
<input type="radio"/>	<input type="radio"/>	0	2734632	U46667.0	Jan 1, 1998 12:30 AM	

<https://submit.ncbi.nlm.nih.gov/>



National Library of Medicine
National Center for Biotechnology Information

Log in

Submission Portal

Submission Portal

Submit to the world's largest public repository of biological and scientific information

Type a few words about the sequence data you are submitting and select an option to learn more. You can also browse submission information below.

What do you want to submit?

Enter a few words about your sequence data.


GenBank


GenBank is the world's largest nucleotide archive containing sequences from all branches of life. The archive is a foundation for medical and biological discovery.




Submitting and updating data

We offer a number of services through which data (including updates) can be submitted to the European Nucleotide Archive (ENA). These technologies provide options appropriate for the scale and frequency of submission, the expertise and capacity of the submitter and the nature of the data to be transferred. The choices below lead users most directly to the appropriate submission route.

 [Submit](#)
[read data](#)

 [Submit](#)
[assembled sequence and/or annotation](#)
(No partial or complete assemblies)

 [Submit](#)
[genome assemblies](#)
(contigs/scaffolds/chromosomes)

 [Email](#)
ENA helpdesk

<https://ddbj.nig.ac.jp/submission>

DDBJ submission portal

Nucleotide

Submission of small-scale nucleotide sequence data with annotation. In case of project data, please use BioProject, MSS, and DRA.

[Create new submission](#)

BioProject

You must obtain BioProject ID and/or locus_tag prefix, before the submission of project data, such as WGS, complete genome, transcriptome project data, DRA, and DTA. The BioProject database collects information about a higher order organization of research projects and its corresponding data. Using BioProject ID make it possible to obtain the same project from various nucleotide sequence databases.

Mass Submission System (MSS)

Please use mass submission system for the submission of following data. WGS, WGS scaffold(s), complete bacterial/eukaryotic genome, HTG, CON, GSS, EST, TSA, and other data includes huge number of sequences.

DDBJ Sequence Read Archive (DRA)

For repository of output data generated by next-generation sequencing machines including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD® System, and others.



Podklady pro GenBank

Povinné položky

- Autoři sekvence
- Kontakt na autory
- Publikace (PubMed)
- Použitá sekvenační technologie
- Použitá metoda assembly
- FASTA nukleotidová sekvence
- Název organismu, taxonomické zařazení
- Metadata
 - izolát, kmen, datum odběru, země původu
- Anotace sekvencí
 - buď vlastní nebo automatická

Protokoly pro zaslání do nukleotidové databáze



- Standard
- Whole Genome Shotgun (WGS)
- Complete Microbial or Eukaryotic Genomes
- ESTs (expressed sequence tags) a GSSs (genome survey sequences)
- High-Throughput Genomic Sequences (HTGs)
- Transcriptome Shotgun Assembly (TSA)
- Third Party Annotation (TPA)
 - záznamy, které upřesňují existující sekvence uložené do databází jinými autory
 - striktní požadavek na přímý experimentální důkaz

Typy standardních anotovaných sekvencí (nucleotide sequence database)

- prokaryotické geny a části genomu
- eukaryotické geny a části genomu
- mRNA sekvence
- rRNA a nebo ITS
- nekódující RNA
- virové sekvence
- transpozony a inzerční sekvence
- mikrosatelity
- pseudogeny
- klonovací vektory
- fylogenetické nebo populační studie (alignments)

Sekvence, které nejsou akceptovány v primárních databázích

- sekvence bez fyzického (biologického) protějšku – např. konsenzní sekvence
- genomové sekvence více exonů bez údajů o sekvencích intronů
- sekvence <200 bp (vyjma patentových)
- sekvence primerů (mohou být zaslány do NCBI's Probe database)
- pouze sekvence proteinů (mohou být zaslány do UniProt/SwissProt)
- sekvence složené z genomové sekvence a mRNA reprezentované jako jedna sekvence

Nezpracovaná zdrojová data z genomových projektů

- BioSample & BioProject mohou obsahovat různé typy archivů

- [Trace Archive](#)

- sekvence získané Sangerovou technikou sekvenování
- struktura složek se *.scf nebo *.abi soubory

```
TOP_DIRECTORY/  
TOP_DIRECTORY/TRACEINFO.txt  
TOP_DIRECTORY/MD5  
TOP_DIRECTORY/README  
TOP_DIRECTORY/traces  
TOP_DIRECTORY/traces/HBBA/  
TOP_DIRECTORY/traces/HBBA/HBBAA1U0001.scf  
TOP_DIRECTORY/traces/HBBA/HBBAA1U0002.scf  
TOP_DIRECTORY/traces/HBBA/HBBAA1U0003.scf
```

- [Sequence Read Archive \(SRA\)](#)

- archiv obsahující alignment sekvencí získaných při 454, IonTorrent, Illumina, SOLiD, Helicos, PacBio nebo Complete Genomics

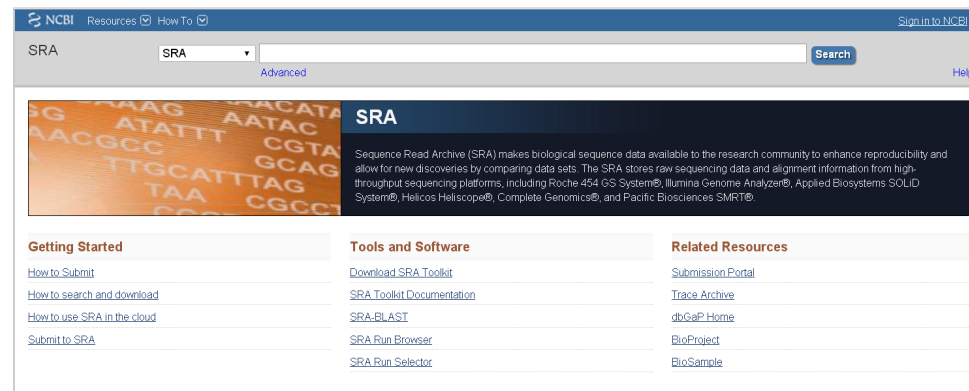
- [The database of Genotypes and Phenotypes \(dbGaP\)](#)

- interakce genotypu a fenotypu člověka

Sequence Read Archive (SRA)

Formát dat a minimální požadavky

- Submission portal
<https://submit.ncbi.nlm.nih.gov/subs/sra/>
- Volitelné nahrání primárních sekvenačních dat
- SRA toolkit
<https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software>
- Minimální požadavek je: primární sekvence (báze) a kvalita = **FASTQ**
- Doporučený formát dat je **BAM** (aligned)
- Další akceptovatelné formáty dat z různých platforem jsou
 - SRF
 - General Fastq
 - SOLiD Fastq
 - Illumina Fastq
 - 454 SFF
 - Ion Torrent SFF
 - PacBio, Nanopore HDF5
 - CompleteGenomics Data Package

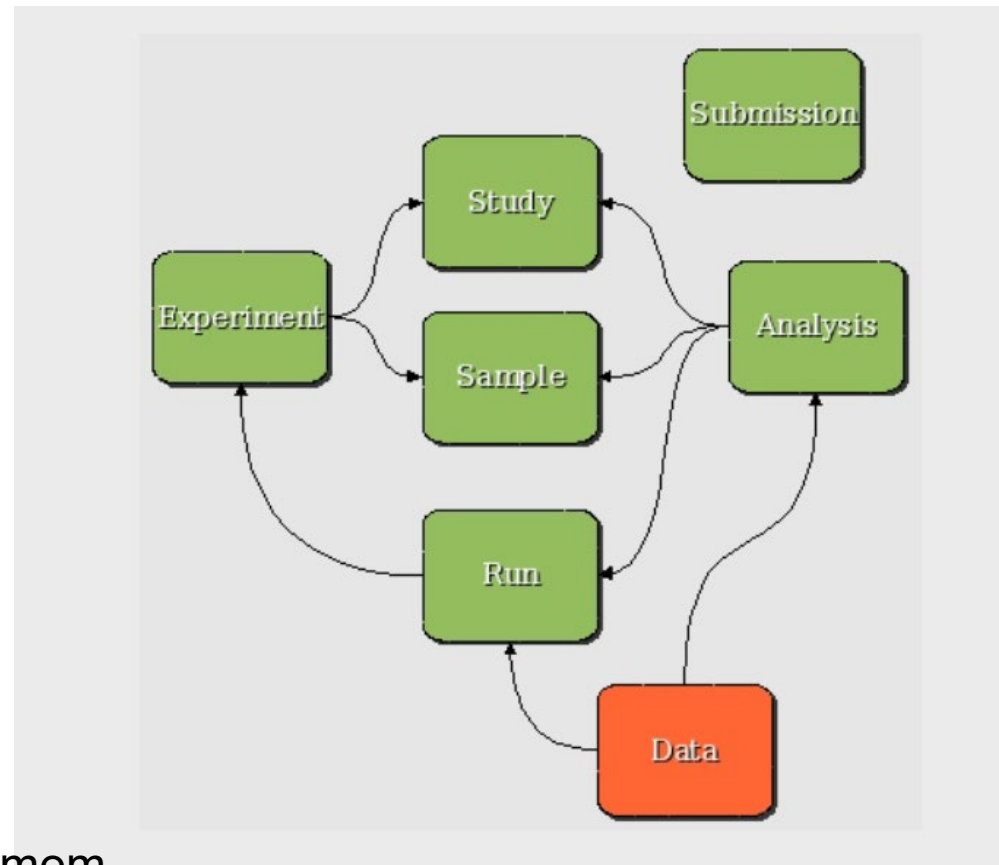


The screenshot shows the NCBI SRA website. At the top, there is a navigation bar with 'NCBI Resources' and 'How To' menus, and a 'Sign in to NCBI' link. Below this is a search bar with 'SRA' entered and a 'Search' button. The main content area features a large image of DNA sequence data (A, T, C, G) and a dark blue box with the text: 'Sequence Read Archive (SRA) makes biological sequence data available to the research community to enhance reproducibility and allow for new discoveries by comparing data sets. The SRA stores raw sequencing data and alignment information from high-throughput sequencing platforms, including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD System®, Helicos Heliscope®, Complete Genomics®, and Pacific Biosciences SMRT®.' Below this, there are three columns of links: 'Getting Started' (How to Submit, How to search and download, How to use SRA in the cloud, Submit to SRA), 'Tools and Software' (Download SRA Toolkit, SRA Toolkit Documentation, SRA-BLAST, SRA Run Browser, SRA Run Selector), and 'Related Resources' (Submission Portal, Trace Archive, dbGap Home, BioProject, BioSample).



Metadata v SRA

- Datové soubory jsou zasílány s metadaty
 - Studie
 - Experiment
 - Vzorek
 - Běh
 - Analýza
 - eticky citlivá data (EGA)



Příklad SRA s mikrobiálním genomem

<https://trace.ncbi.nlm.nih.gov/Traces/sra/?run=SRR9600155>

Whole Genome Shotgun (WGS)

- WGS sekvenační projekty jsou celé genomy nebo chromozomy sekvenované strategií celogenomového shotgun sekvenování
- Části WGS projektu jsou kontigy, které nesmí obsahovat mezery
- WGS projekty mohou být anotovány, může být zvolena automatická anotace s NCBI pipeline
- Volitelně - soubor [AGP](#) ukazuje, jak jsou kontigy oddělené mezerami uspořádány na chromozomu
- Zasílají se přes Genome submission portal

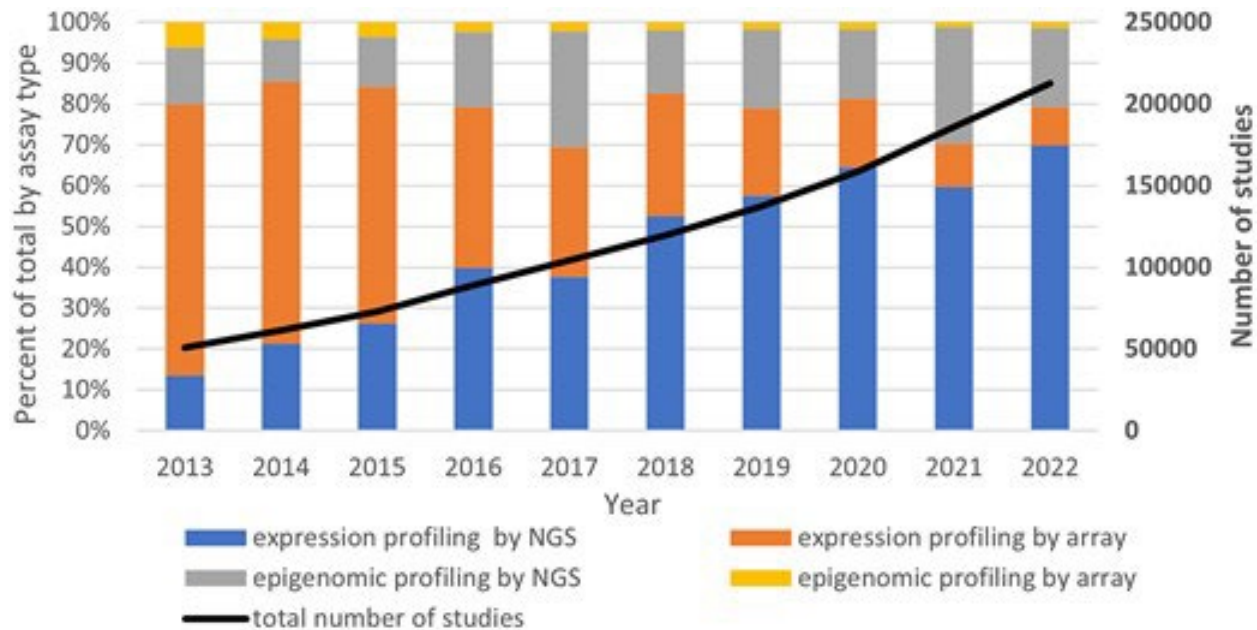
Metagenomy

- Sekvenační projekty analýzy společenstev z určitých ekologických zdrojů nezávislé na kultivaci slouží pro studium
 - genetické diversity, struktury populací, ekologické úlohy
 - metabolických funkcí
 - stanovení kompletních genomů nekultivovatelných organismů
 - izolaci nových mikroorganismů (genetických zdrojů) z prostředí
- **Sekvence jsou vzájemně propojené v rámci BioProject**
- **Bývá vyžadováno nahrání do SRA**
- Metagenomové projekty mohou představovat
 - Neanotované sekvence s převahou informačních sekvencí rRNA
 - Celkové metagenomové projekty sestavené do kontigů
 - obsahují částečné genomy z taxonomicky různých skupin
 - MG-RAST - anotace metagenomu na serveru
- Mothur a QIIME2 obsahují nástroje pro přípravu dat

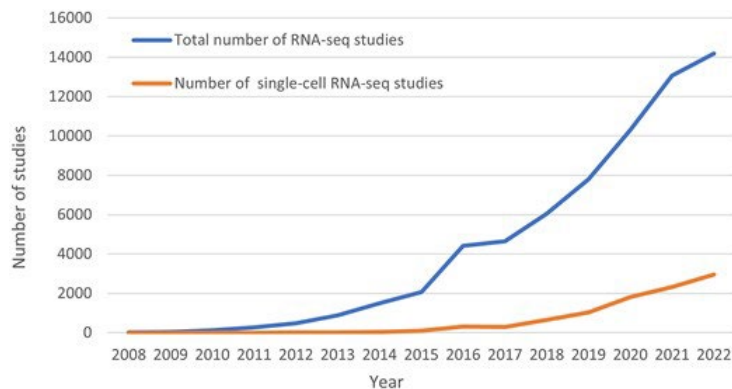
NCBI Gene Expression Omnibus (GEO)

- mezinárodní veřejné úložiště v NCBI
 - <https://www.ncbi.nlm.nih.gov/geo>
- archivuje zpracované datové soubory a metadata o
 - genové expresi
 - epigenomice
- generované technologiemi
 - sekvenování nové generace
 - Microarray
- nabízí webové nástroje, pro analýzu a vizualizaci diferenciální genové exprese

NCBI Gene Expression Omnibus (GEO)

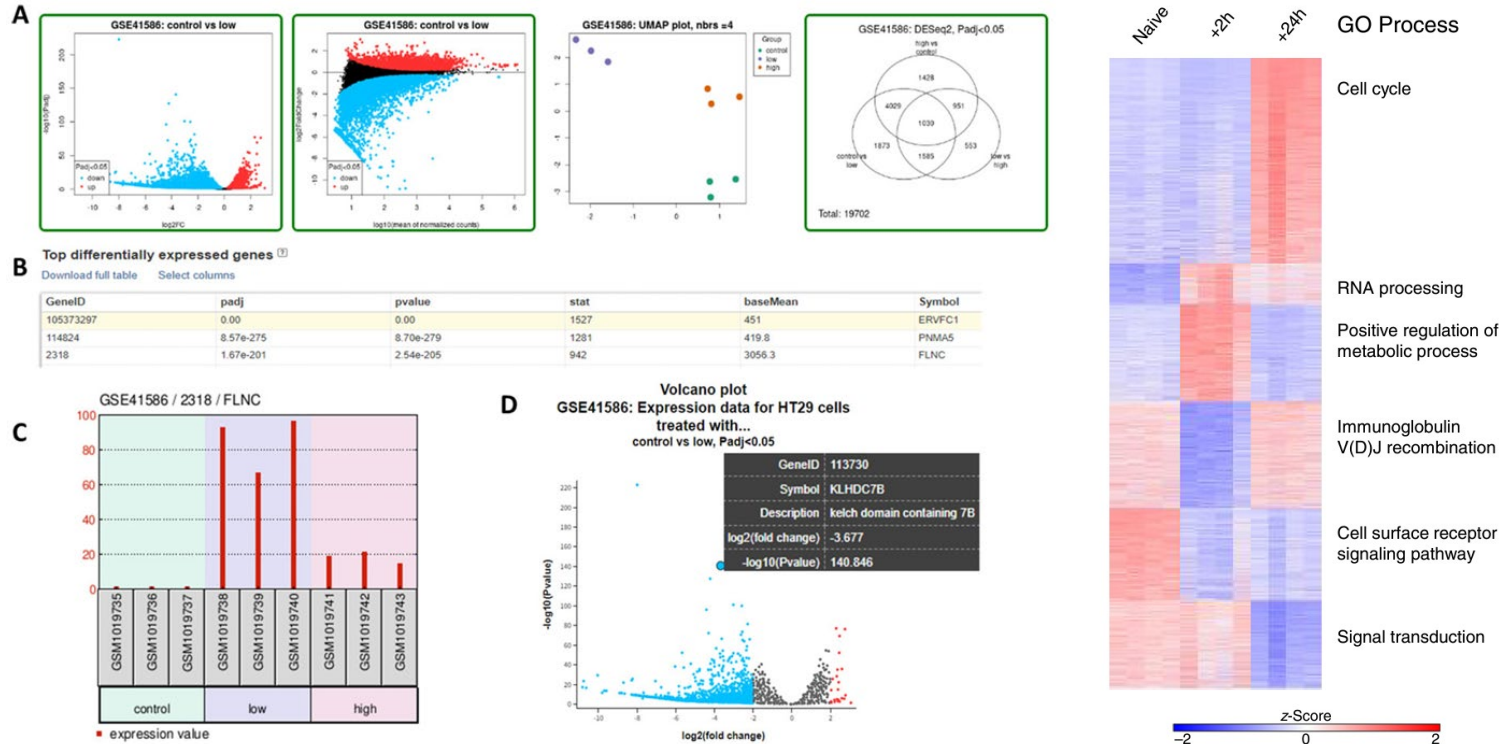


Růst a trendy datových typů za poslední desetiletí v databázi GEO



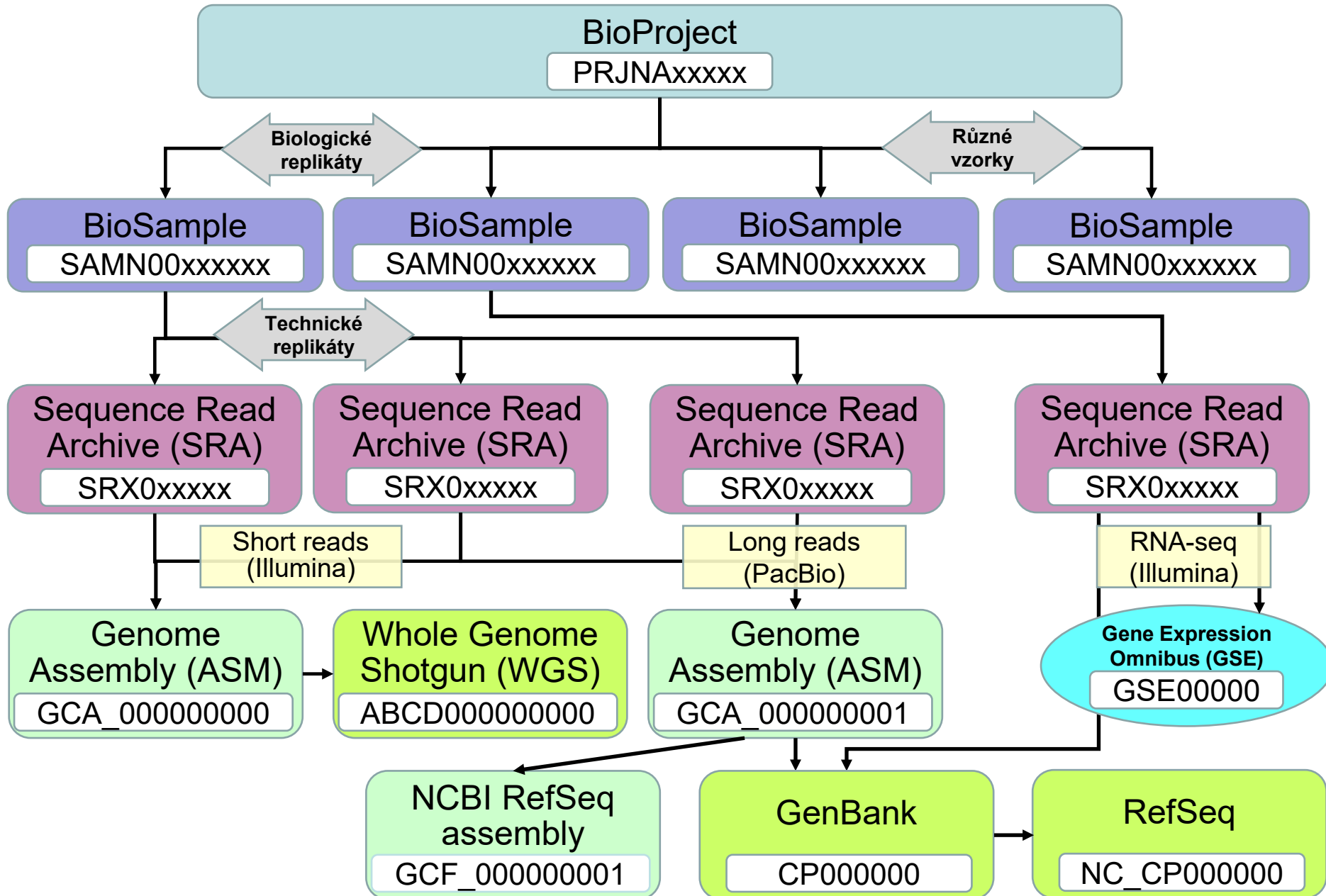
Z toho počet studií RNA-seq a jednobuněčné RNA-seq v GEO

NCBI Gene Expression Omnibus (GEO)



Příklady vizualizace dat z databáze GEO nástrojem GEO2R. Poskytuje interaktivní grafy (Umístěním kurzoru na bod se zobrazí jeho GeneID, Symbol, Popis, \log_2 (násobná změna) a $-\log_{10}$ (P-hodnota), tabulky nejlepších odlišně exprimovaných genů se statistikou („Top“ diferenciálně exprimované geny)

Struktura databází (BioSample & BioProject)

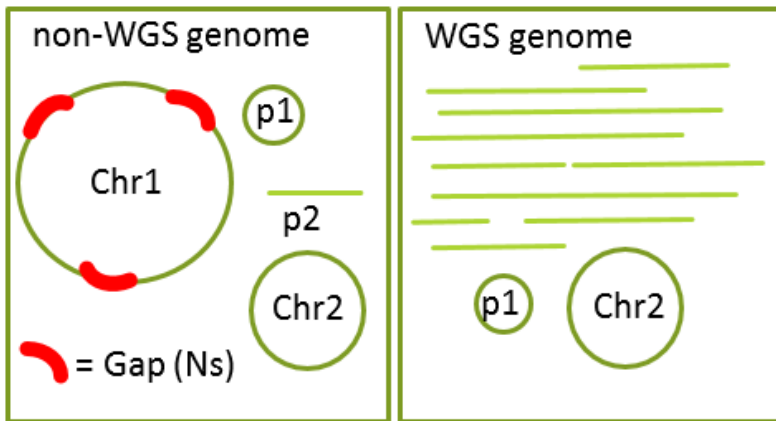


Genome submission portal

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

Prokaryotic and Eukaryotic Genomes Submission Guide

Both WGS and non-WGS genomes, including gapless complete bacterial chromosomes, can be submitted via the Submission Portal. You will be asked to choose whether the genome being submitted is considered WGS or not. The differences for GenBank purposes are:



non-WGS

- Each chromosome is in a single sequence and there are no extra sequences
- Each sequence in the genome must be assigned to a chromosome or plasmid or organelle
- Plasmids and organelles can still be in multiple pieces.

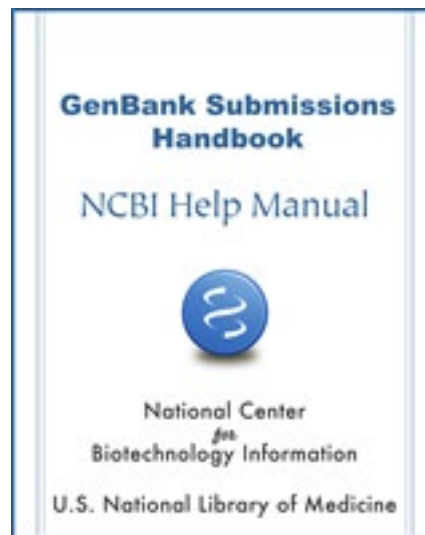
WGS

Genome Resources

- [About WGS](#)
- [WGS Browser](#)
- [Genome Submission Guide](#)
- [Genome Submission Portal](#)
- [Update Genome Records](#)
- [FAQ](#)
- [tbl2asn](#)
- [Create Submission Template](#)
- [Eukaryotic Annotation Guide](#)
- [Prokaryotic Annotation Guide](#)
- [Annotation Example Files](#)
- [Discrepancy Report](#)
- [NCBI Prokaryotic Genome Annotation Pipeline](#)
- [AGP Format](#)
- [Complex Assembly Submission Guide](#)
- [Metagenome Submission Guide](#)
- [BioProject](#)


Postup zaslání GenBank Standardního typu

- <http://www.ncbi.nlm.nih.gov/books/NBK51157/>
The GenBank Submissions Handbook



Submit new sequences to GenBank

What type of sequence data do you have?

- SARS-CoV-2 
- Ribosomal RNA (rRNA) or rRNA-ITS
- Metazoan (multicellular animal) COX1
- Influenza virus
- Norovirus
- Dengue virus
- Eukaryotic and Prokaryotic Genomes (WGS or Complete)
- Transcriptome Shotgun Assembly (TSA)
- Unassembled sequence reads (SRA)

- Sequence data not listed above (through BankIt): mRNA, genomic DNA, organelle, ncRNA, plasmids, other viruses, phages, synthetic constructs

Start

Další požadavky na zaslání sekvence

- Informace o datu zveřejnění
- Informace o relevantních publikacích
- Popis zdroje sekvence
- Vlastní sekvence
 - typ a tvar molekuly
 - anotace vlastností sekvence

Popis zdroje sekvence 1

- **organism**
nezkrácené vědecké jméno
Příklad: [organism=Drosophila melanogaster]
- **lineage**
taxonomické zařazení organismu (dle NCBI taxonomy database)
<http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Root>
- **molecule**
ve tvaru "DNA" nebo "RNA".
Příklad : [molecule=DNA]
- **moltype**
může nabývat následujících hodnot
Příklad : [moltype=Genomic DNA]
 - Genomic DNA
 - Genomic RNA
 - Precursor RNA
 - mRNA [cDNA]
 - Ribosomal RNA
 - Transfer RNA
 - Small nuclear RNA
 - Small cytoplasmic RNA
 - Other-Genetic
 - cRNA
 - Small nucleolar RNA
- **topology**

Popis zdroje sekvence 2

- **location**
může nabývat následujících hodnot
Příklad: [location=mitochondrion]
 - genomic
 - chloroplast
 - kinetoplast
 - mitochondrion
 - plastid
 - macronuclear
 - extrachromosomal
 - plasmid
 - cyanelle
 - proviral
 - virion
 - nucleomorph
 - apicoplast
 - leucoplast
 - proplastid
 - endogenous-virus
 - hydrogenosome
- **Genetic code**
<http://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi?mode=c>

Popis zdroje sekvence 3

Další popisovače ke zdroji sekvence

- acronym
- anamorph
- authority
- biotype
- biovar
- breed
- cell-line
- cell-type
- chemovar
- chromosome
- clone
- clone-lib
- collected-by
- common
- country
- cultivar
- dev-stage
- ecotype
- endogenous-virus-name
- forma
- forma-specialis
- fwd-pcr-primer-name
- fwd-pcr-primer-seq
- genotype
- group
- haplotype
- identified-by
- isolate
- isolation-source
- lab-host
- lat-lon
- map
- note
- pathovar
- plasmid-name
- plastid-name
- pop-variant
- rev-pcr-primer-name
- rev-pcr-primer-seq
- segment
- serogroup
- serotype
- serovar
- sex
- specific-host
- specimen-voucher
- strain
- sub-species
- subclone
- subgroup
- substrain
- subtype
- synonym
- teleomorph
- tissue-lib
- tissue-type
- type
- variety

Formát sekvence

- Sekvence nukleové kyseliny a kódovaných proteinů připravené ve formátu FASTA

Nucleotide Sequence:

```
>ABC-1 [organism=Saccharomyces cerevisiae][strain=ABC][clone=1]
ATTGCGTTATGGAAATTCGAAACTGCCAAATACTATGTCACCATCATTGA
TGCACCTGGACACAGAGATTTTCATCAAGAACATGATCACTGGTACTT
```

Protein Sequences:

```
>4E-I [gene=eIF4E] [protein=eukaryotic initiation factor 4E-I]
MQSDFHRMKNFANPKSMFKTSAPSTEQGRPEPPTSAAAPAEAKDVKPKEDPQETGEPAGN ...
>4E-II [gene=eIF4E] [protein=eukaryotic initiation factor 4E-II]
MVVLETEKTSAPSTEQGRPEPPTSAAAPAEAKDVKPKEDPQETGEPAGNTATTTAPAGDD ...
```

Přerušená sekvence

```
>m_gagei [organism=Mansonia gagei] Mansonia gagei NADH dehydrogenase ...
ATGGAGCATACATATCAATATTCATGGATCATAACGTTTGTGCCACTTCCAATTCCTATTTTAATAGGAA
TTGGACTCCTACTTTTTCCGACGGCAACAAAAAATCTTCGTCGTATGTGGGCTCTTCCCAATATTTTATT
GTTAAGTATAGTTATGATTTTTTCGGTCGATCTGTCCATTCAGCAAATAAATAAAAGTTCTATCTATCAA
TATGTATGGTCTTGGACCATCAATAATGATTTTTCTTTCGAGTTTGGCTACTTTATTGATTCGCTTACCT
>?200 ← Délka přerušení
GGTATAATAACAGTATTATTAGGGGCTACTTTAGCTCTTGC
TCAAAAAGATATTAAGAGGGGTTTAGCCTATTCTACAATGTCCCAACTGGGTTATATGATGTTAGCTCTA
GGTATGGGGTCTTATCGAGCCGCTTTATTTCAATTTGATTACTCATGCTTATTTCGAAGGCATTGTTGTTTT
TAGGATCCGGATCCGTTATTCATTCCATGGAAGCTATTGTTGGATATTCTCCAGATAAAAGCCAGAATAT
GGTTTTTATGGGCGGTTTAAGAAAGCATGTGCCAATTACACAAATTGCTTTTTTTAGTGGGTACACTTTCT
CTTTGTGGTATTCACCCCTTGCTTGTTTTTTGGTCCAAAGATGAAATTCCTTAGTGACAGCTGGTTGT
>?unk100 ← Přerušení neznámé délky
TCAATAAAACTATGGGGTAAAGAAGAACAAAAATAATTAACAGAAATTTTCGTTTATCTCCTTTATTAA
TATTAACGATGAATAATAATGAGAAGCCATATAGAATTGGTGATAATGTAAAAAAGGGGCTCTTATTAC
TATTACGAGTTTTGGCTACAAGAAGGCTTTTTCTTATCCTCATGAATCGGATAATACTATGCTATTTCCCT
ATGCTTATATTGGCTCTATTTACTTTTTTTGTTGGAGCCATAGCAATTCCTTTTAATCAAGAAGGACTAC
ATTTGGATATATTATCCAAATTATTA ACTCCATCTATAAATCTTTTACATCAAATTC AAATGATTTTGA
GGATTGGTATCAATTTTTAACA AATGCAACTCTTTCAGTGAGTATAGCCTGTTTCGGAATATTTACAGCA
TTCTTTTTATATAAGCCTTTTTTATTCATCTTTACAAAATTTGA ACTTACTAAATTTATTTTCGAAAGGG
GTCCTAAAAGAATTTTTTTGGATAAAATAATACTTGATATACGATTGGTCATATAATCGTGGTTACAT
```

Sekvenční příložen

- Fasta+GAP

```
>ABC-1 [organism=Saccharomyces cerevisiae][strain=ABC][clone=1]
---ATTGCGTTATGGAAATTCGAAACTGCCAAATACTATGTCACCATCAT
TGATGCACCTGGACACAGAGATTTTCATCAAGAACATGATCACTGGTACTT
>ABC-2 [organism=Saccharomyces cerevisiae][strain=ABC][clone=2]
GATATTGCTTTATGGAAATTCGAAACTGCCAAATACTATGTCACCATCAT
TGATGCACCTGGACACAGAAATTTTCATCAAGAACATGATCACTGGTACTT
>ABC-3 [organism=Saccharomyces cerevisiae][strain=ABC][clone=3]
---ATTGCTTTATGGAAATTCGAAACTGCCAAATACTATGTTA-----
TGATGCACCTGGACACAGAGATTTTCATCAAAAACATGATCACTGGTACTT
```

- PHYLIP

```
3 100
ABC-1 ---ATTGCGT TATGGAAATT CGAAACTGCC AAATACTATG TCACCATCAT
ABC-2 GATATTGCTT TATGGAAATT CGAAACTGCC AAATACTATG TCACCATCAT
ABC-3 ---ATTGCTT TATGGAAATT CGAAACTGCC AAATACTATG TTA-----

TGATGCACCT GGACACAGAG ATTTTCATCAA GAACATGATC ACTGGTACTT
TGATGCACCT GGACACAGAA ATTTTCATCAA GAACATGATC ACTGGTACTT
TGATGCACCT GGACACAGAG ATTTTCATCAA AAACATGATC ACTGGTACTT
```

```
>[organism=Saccharomyces cerevisiae][strain=ABC][clone=1]
>[organism=Saccharomyces cerevisiae][strain=ABC][clone=2]
>[organism=Saccharomyces cerevisiae][strain=ABC][clone=3]
```

Anotace vlastní sekvence

- Kódované proteiny
 - CDS
interval
nekompletnost na N- nebo C- konci
 - gene
interval odpovídající CDS u experimentálně prokázaných genů
 - mRNA
interval obsahující 5'-UTR a 3'-UTR
- Kódované strukturní RNA

Příklady některých dalších modifikací deskriptorů

- Title
 - Informace vyskytující se v databázi v DEFINITION LINE
- Comment
 - Poznámka k různým vlastnostem
- Technique
 - Umožňuje výběr techniky použité pro vytvoření nebo experimentální evidenci vlastností sekvence

Přehled deskriptorů pro popis vlastností sekvence

(<http://www.ncbi.nlm.nih.gov/BankIt/help.html>)

- attenuator
- C-region
- CAAT_signal
- CDS
- conflict
- D-loop
- D-segment
- enhancer
- exon
- gap
- GC_signal
- gene
- iDNA
- intron
- J_segment
- LTR
- mat_peptide
- misc_binding
- misc_difference
- misc_feature
- misc_recomb
- misc_RNA
- misc_signal
- misc_structure
- modified_base
- mRNA
- N_region
- old_sequence
- operon
- oriT
- polyA_signal
- polyA_site
- precursor_RNA
- prim_transcript
- primer_bind
- promoter
- protein_bind
- RBS
- repeat_region
- repeat_unit
- rep_origin
- rRNA
- S_region
- satellite
- scRNA
- sig_peptide
- snRNA
- snoRNA
- source
- stem_loop
- STS
- TATA_signal
- terminator
- transit_peptide
- tRNA
- unsure
- V_region
- V_segment
- variation
- 3'clip
- 3'UTR
- 5'clip
- 5'UTR

Příklady sekvencí

Sekvence mRNA nebo cDNA

- Kódující oblasti včetně iniciačního a terminačního kodonu
- Název proteinu
- Název genu
- Sekvence proteinu

Homo sapiens prolidase (PEPD) mRNA, complete cds.

FEATURES	Location/Qualifiers
source	1..1888 /organism="Homo sapiens" /chromosome="19" /map="19q12-q13.2" /cell_type="fibroblasts"
mRNA	1..1888 /gene="PEPD"
gene	1..1888 /gene="PEPD"
CDS	17..1498 /gene="PEPD" /EC_number="3.4.13.9" /note="imidodipeptidase" /product="prolidase"

Sekvence prokaryotického genu

- Kódující intervaly
- Název proteinu
- Název genu, je-li známý
- Aminokyselinová sekvence

`Escherichia coli RecA protein (recA) gene, complete cds.`

<code>FEATURES</code>	<code>Location/Qualifiers</code>
<code>source</code>	<code>1..3300</code> <code>/organism="Escherichia coli"</code> <code>/strain="K-12"</code>
<code>gene</code>	<code>783..1961</code> <code>/gene="recA"</code>
<code>CDS</code>	<code>783..1961</code> <code>/gene="recA"</code> <code>/function="DNA repair protein"</code> <code>/product="RecA protein"</code>

Ribosomální RNA a vnitřní přepisované mezeríky

- Názvy jakékoli strukturní RNA (např. tRNA-Ile, 16S ribosomal RNA)
- Názvy mezeríkových oblastí (např., internal transcribed spacer 1, 16S/23S intergenic spacer)
- Nukleotidové pozice

`Saccharomyces cerevisiae 18S ribosomal RNA gene, partial sequence; internal transcribed spacer 1, 5.8S ribosomal RNA gene and internal transcribed spacer 2, complete sequence; and 28S ribosomal RNA gene, partial sequence.`

FEATURES	Location/Qualifiers
source	1..540 /organism="Saccharomyces cerevisiae" /strain="UMD 334"
rRNA	<1..5 /product="18S ribosomal RNA"
misc_RNA	6..178 /product="internal transcribed spacer 1 "
rRNA	179..377 /product="5.8S ribosomal RNA"
misc_RNA	378..519 /product="internal transcribed spacer 2"
rRNA	520..>540 /product="28S ribosomal RNA"

Oblast promotoru

- Název proteinu nebo genu, ke kterému patří promotor a jeho 5' a 3' obklopující sekvence
- Intervaly přepisovaných a kódujících sekvencí, pokud jsou přítomné

Homo sapiens enhancer-binding protein 2 (EBP2) gene, promoter region and partial cds.

FEATURES	Location/Qualifiers
source	1..3061 /organism="Homo sapiens" /chromosome="15" /map="15q13" /cell_line="H441" /tissue_type="lung"
gene	1..>3061 /gene="EBP2"
promoter	1..2947 /gene="EBP2"
TATA_signal	2918..2923 /gene="EBP2"
mRNA	2948..>3061 /gene="EBP2" /product="enhancer-binding protein 2"
5' UTR	2948..3010 /gene="EBP2"
CDS	3011..>3061 /gene="EBP2" /product="enhancer-binding protein 2"

Transpozon nebo inzerční sekvence

Specifické jméno elementu

- Nukleotidové pozice
- Jména a intervaly kódovaných genových produktů, pokud jsou přítomny (např., transposase)
- Pozice a intervaly dalších vlastností (např. LTRs, repeat regions)

**Bacillus subtilis transposon BLT transposase (tnpA) gene,
complete cds**

```
FEATURES             Location/Qualifiers
    source             1..1221
                       /organism="Bacillus subtilis"
                       /strain="RS2"
    source             21..1127
                       /organism="Bacillus subtilis"
                       /strain="RS2"
                       /transposon="BLT"
    repeat_region      21..61
                       /rpt_type=inverted
    gene               128..1034
                       /gene="tnpA"
    CDS                128..1034
                       /gene="tnpA"
                       /product="transposase"
    repeat_region      1085..1127
                       /rpt_type=inverted
```

Oblasti repeticí

- Intervaly repetitivních sekvencí
- Rodina repeticí (např., Alu, Mer)
- Typ repetice (tandem, inverted, flanking, terminal, direct, dispersed, or other)
- Jednotka repetice (repeat unit) popis intervalů, jestliže sekvence obsahuje více než jednu repetici

Homo sapiens repeat regions

FEATURES	Location/Qualifiers
source	1..2050 /organism="Homo sapiens" /chromosome="6" /map="6q25"
repeat_region	8..126 /rpt_type=dispersed /rpt_family="B2"
repeat_region	197..344 /rpt_type="direct" /rpt_unit="197..220"
repeat_region	389..673 /rpt_family="AluSx" /rpt_type=dispersed
repeat_region	847..876 /note="microsatellite BT21" /rpt_type="tandem" /rpt_unit="ca"
repeat_region	1000..2000 /rpt_family="human endogeneous retrovirus K-10"

Klonovací vektor

- Jedinečné jméno vektoru
- Kódující intervaly, jména genů a proteinů

Cloning vector pRB223, complete sequence

FEATURES	Location/Qualifiers
source	1..4361 /organism="Cloning vector pRB223"
gene	86..1276 /gene="tet"
CDS	86..1276 /gene="tet" /product="tetracycline resistance protein"
RBS	1905..1909 /note="Shine-Dalgarno sequence"
rep_origin	2535
gene	complement(3293..4194) /gene="bla"
CDS	complement(3293..4153) /gene="bla" /product="beta-lactamase"
misc_feature	4069..4125 /note="multiple cloning site"
RBS	complement(4161..4165) /gene="bla" /note="Shine-Dalgarno sequence"
promoter	complement(4188..4194) /gene="bla"

