

MUNI

# Data v mikrobiologii

CORE042: Data – odpověď na základní otázku života, vesmíru a vůbec...

6. přednáška

**Stanislava Bezdíček Králová** <[kralova.s@sci.muni.cz](mailto:kralova.s@sci.muni.cz)> a kol.

Externista – Přírodovědecká fakulta MU

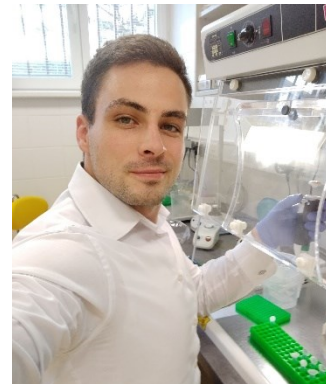
# Přednášející a spoluautoři



- **Stanislava Bezdíček Králová**
  - University of Vienna, Center of Microbial Ecology and System Science
  - Přírodovědecká fakulta MU, externista
  - **Life science projekty, FAIR Data v mikrobiologii**



- **Václav Šeda**
  - CEITEC
  - Mikroprostředí imunitních buněk
  - **FAIR data v klinickém a univerzitním výzkumu**



- **Matěj Bezdíček**
  - Fakultní nemocnice Brno
  - Interní hematologická a onkologická klinika
  - **FAIR data v klinické mikrobiologii**

# Osnova

- Life cycle
- Historie vs. současnost
- Plánování
- Sběr/získávání dat
- Zpracování dat
- Analýza dat
- Uchování dat
- Poskytování dat
- Opakované využití dat



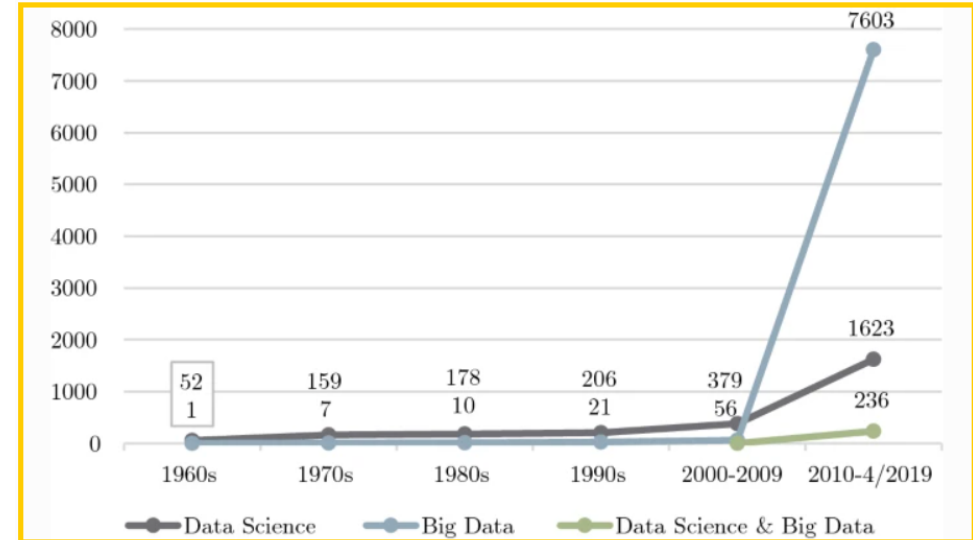
# Historie

## V minulosti

- méně dat
- menší výpočetní náročnost
- jednodušší na uchování, organizaci, vyhodnocení

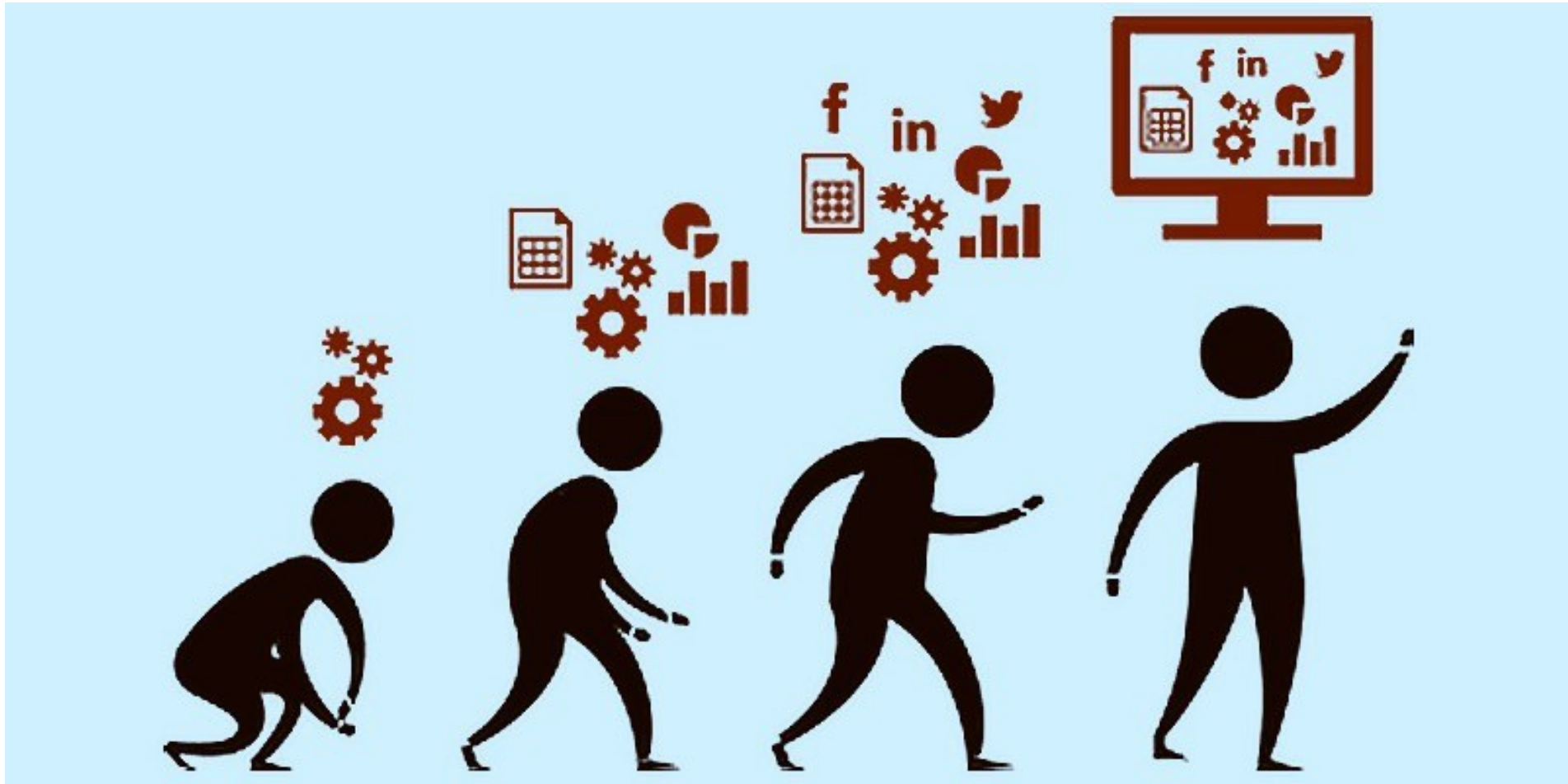
## Současnost

- změna instrumentace
- změna charakteru dat (pozorování vs. sekvenace)
- enormní nárůst množství dat – digitální forma
- vyšší nároky na rutinní procesování – zisk, organizace dat, jejich sdílení, uchovávání, procesování (úpravy), analýzy



Evolutionary trend in the number of publications covering data science and big data

*Raban & Gordon, Scientometrics (2020)*





# Plánování



## Co?

- stanovení cíle

## Jak?

- určení metodologického postupu  
(*máme vše co potřebujeme?*)
- určení znalostního základu  
(*dokážeme uskutečnit všechny kroky?*)
- dohodnutí spolupráce
- rozložení finančních prostředků

## Kdo?

- jeden výzkumník?
- tým?
- studenti?
- kdo to zaplatí?!

## Kde?

- získat vzorek + logistika
- máme potřebnou instrumentaci?
- dílčí části u spolupracujícího pracoviště?  
(tuzemsko/zahraníčí)

## Kdy?

- časové plánování je extrémně důležité
- kdy odeberu vzorky
- stihnu je pak zpracovat?
- kolik vzorků zpracuji?
- potřebuji více odběrů v různých časech?
- dílčí kroky vs. celkový projekt

# Plánování – specifika klinického výzkumu



## Etická komise

- Lidská DNA ve vzorcích = etický problém  
→ nutné schválení
- Ochrana osobních údajů
- Informovaný souhlas

## Původ vzorků?

- ambulantní vs. hospitalizovaný pacient
- transport
- krátkodobé uchování
- dlouhodobé uchování
- jeden výzkumník?
- tým?
- studenti?
- kdo to zaplatí?!

## Kdy?

- časové plánování je extrémně důležité
- jak často budeme odebírat vzorky?
- hledám něco typické nebo spíš raritní??  
→ někdy extrémně dlouhá doba do vytvoření dostačujícího souboru vzorků

## Kdo?

- interpretace: často nutný klinik
- často nutné zapojení více pracovišť

## Financování

- nemocnice
- granty → klinický vs základní výzkum



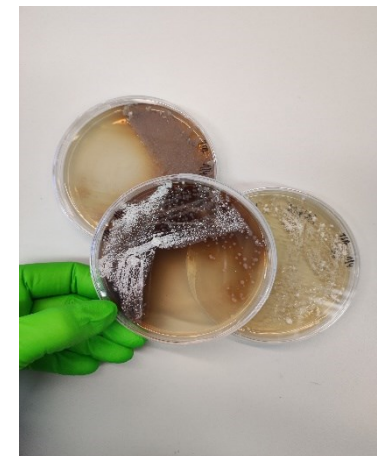
# Plánování – příklady

## Terénní výzkum (Antarktida)

- chystání materiálu (6 měsíců před expedicí)
- chystání citlivního materiálu (1 měsíc před expedicí)
- jak budou vzorky cestovat?
- jaké vybavení budu mít k dispozici?
- jak dokážu zpracovat vzorky?

## Přílet na Antarktidu:

1. příprava kultivačních médií
2. sterilizace materiálu po cestě
3. odběry **na začátku** expedice (kultivace 6–8 týdnů)
4. zpracovávání vzorků
5. odběry **na konci** expedice (přenos do ČR)
  - ✓ nesmí roztát
  - ✓ nesmí se zahřát
  - ✓ nesmí se střídat teploty
  - ✓ kolik gramů
  - ✓ kolik replikátů
  - ✓ původ vzorků





# Když plánování experimentu nevyjde – příklad z klinické praxe

## Studium invazivních mykóz

- Zahájení nového projektu současně se sběrem vzorků
- Spolupráce s klinikou v rámci nemocnice
- Za rok sběru vzorků nasbírám 2 potvrzené a 3 suspektní vzorky na invazivní mykózu
  
- výsledky → ***z výsledků u 5 vzorků neudělám relevantní závěry, nic nepublikuji***

## Studium invazivních mykóz

- **Zahájení spolupráce s dalšími centry a po dlouhodobém bankování vzorků**
- Zahájení projektu s již získaným větším množstvím vzorků
  
- výsledky → ***signifikantní počet vzorků k vyšetření umožňuje výsledky hodnotit, korelovat a dělat statisticky významné závěry, publikaci práce založené na takovém souboru nic nebrání***



# Když plánování experimentu nevyjde – příklad z praxe

## Myší experiment – nepromyšlené

- ráno: podání sulfoquinolózy žrádla (10mg/l)
- kontrolní skupina – bez přídavku
- po 12 hodinách odběr stolice
- sekvenování
  
- výsledky → **žádný rozdíl** kontrola vs. testovaná skupina

## Myší experiment – kvalitní plán

- **večer** podání sulfoquinolózy žrádla (10mg/l)
- kontrolní skupina – bez přídavku
- po **3, 6, 12 a 24** hodinách odběr stolice
- sekvenování každé vzorky
  
- výsledky → **rozdíl** mezi skupinami **pouze v 6 h**





# Získávání dat



## Terén

- pozorování + zápis
- měření + zápis

## Laboratoř

- pozorování + zápis
- měření + zápis
- generace dat přístrojem → hrubá data



***experimentálně (původce = vy)***

## Sdílená data

- použití dat v databázích s open access
- použití poskytnutých dat
- „share and re-use“

## Placená data

- placené databáze



***původce je někdo jiný***



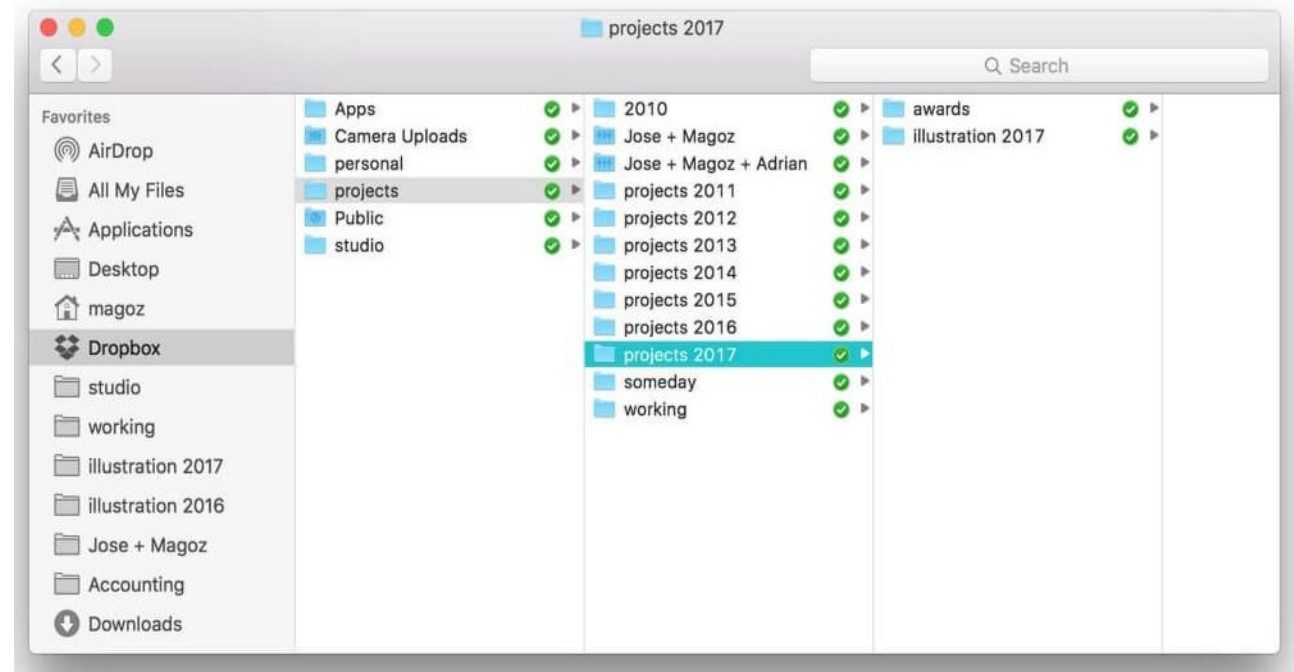
# Organizace dat

## Složky → podsložky → soubory

- konzistentní pojmenování (název projektu, název analýzy, datum provedení analýzy)
- hierarchický systém
- oddělení probíhajících a uzavřených projektů
- najít si vlastní systém který jste schopní **dlouhodobě sledovat a udržet**

## Pojmenování souborů

- konzistentně
- datумы v určitém sledu
- interpunkce (\_ / . mezery???)
- číselné pořadí na začátku / na konci





# Organizace dat

## Zlatý standard – tabulky

- jasně
- stručně
- jedna poznámka jedno místo
- stejný formát v daném sloupci
- jeden nadpis na sloupec
- žádné slučování buněk
- žádné prázdné buňky
- kde to jde dělat vybírací seznamy  
(*víc lidí pracuje se stejnou tabulkou*)
- žádné výpočty v raw data souboru



	A	B	C	D	E	F
1	Region	Retailer Type	Qtr1	Qtr2	Qtr3	Qtr4
2	Mid West					
3		Food & Staples	19,47,100	26,93,000	19,50,300	19,11,400
4		Multiline	9,77,200	7,29,700	3,67,800	3,22,600
5	North East					
6		Specialty	14,55,100	14,22,700	4,98,200	17,86,900
7	South					
8		Food & Staples	16,91,700	22,39,800	18,30,000	17,21,800
9		Multiline	18,37,500	18,80,500	17,20,500	24,89,200
10		Specialty	11,55,100	21,00,200	16,52,100	9,07,700
11	West					
12		Food & Staples	12,50,900	23,68,200	15,55,100	14,55,500
13		Multiline	13,74,300	7,93,200	12,54,900	15,17,400

**BAD DESIGN**

	A	B	C	D	E
1	Region	Qtr1	Qtr2	Qtr3	Qtr4
2	Mid West	2924300	3422700	2318100	2234000
3	North East	1455100	1422700	498200	1786900
4	South	4684000	6220500	5202600	5118700
5	West	2625200	3161400	2810000	2972900

**BAD DESIGN**

	A	B	C	D	E	F	G	H	I
1	Region	Qtr1	Qtr2	Qtr3	Qtr4	Region	Quarter	Sales	
2	Mid West	2924300	3422700	2318100	2234000	Mid West	Qtr1	2924300	
3	North East	1455100	1422700	498200	1786900	Mid West	Qtr2	3422700	
4	South	4684000	6220500	5202600	5118700	Mid West	Qtr3	2318100	
5	West	2625200	3161400	2810000	2972900	Mid West	Qtr4	2234000	
6						North East	Qtr1	1455100	
7						North East	Qtr2	1422700	
8						North East	Qtr3	498200	
9						North East	Qtr4	1786900	
10						South	Qtr1	4684000	
11						South	Qtr2	6220500	
12						South	Qtr3	5202600	
13						South	Qtr4	5118700	
14						West	Qtr1	2625200	
15						West	Qtr2	3161400	
16						West	Qtr3	2810000	
17						West	Qtr4	2972900	

# Organizace dat a metadat

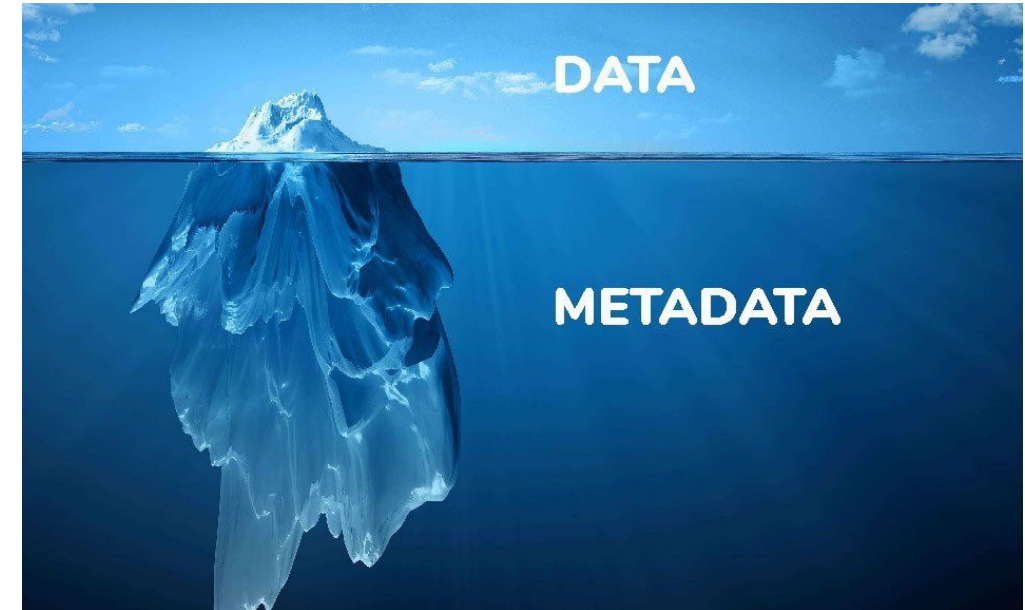


## Doplňkové informace

- mohou být uloženy separátně od metadat nebo dat
  - dotazníky (pacienti, probandi)
  - **laboratorní deníky**
  - **metadata**
  - reporty
  - publikace

## Rozšiřující informace

- přesné popisky míst z tabulky
- přesný popis odběru vzorků z míst v tabulce







# Doplňkové informace → nutnost?!

## Laboratorní deníky

- novinka grantových agentur i některých pracovišť (zejména zahraničí)
- každý člen týmu vede laboratorní deník – přesné záznamy
- skladování i po skončení grantu
- varianty: papírová i elektronická



- naučit se zaznamenávat práci v laboratoři
- případné změny = lehčí implementace
- nutnost u projektů i grantů, i když to nevyžadují

## Příklad: Rakousko – přísná pravidla

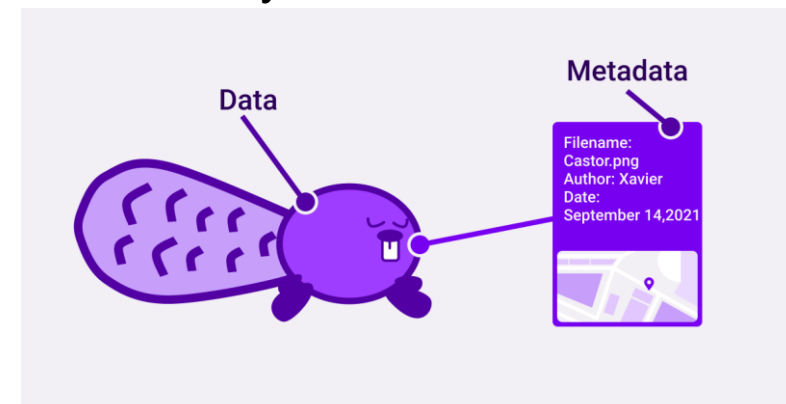
- denně podrobné zápisy o průběhu experimentů
- žádné škrtní
- žádné trhání listu
- elektronická verze musí být opatřena časovým zámkem

# Metadata v biologii



## Metadata = informace o nasbíraných datech

- **Fotky:** jak-kde-kdy vznikli, čím byli pořízeny, barevný profil, čočky, zoom...
- **Vzorek:** kde se odebral, GPS, kolik, typ zdrojového materiálu, jak se skladoval, jak se transportoval, jak se zpracoval
- **Bakteriální izolát:** Z jaké vzorky pochází, za jaké teploty roste, na jakých mikrobiologických půdách roste, do jaké BSL kategorie spadá
- **Experimentální data:** design jednotlivých experimentů, teploty, pH, živiny, obsah vody, salinita, objemy, hmotnosti, trvání, jakékoliv změny



# Metadata v biologii



## Sekvenační data:

### Links from BioProject

#### Flavobacterium flabelliformis sp. nov. P4023 T

Identifiers BioSample: SAMN18262075; Sample name: Flavobacterium flabelliformis sp. nov. P4023T

Organism [Flavobacterium flabelliforme](#)

cellular organisms; Bacteria; FCB group; Bacteroidetes/Chlorobi group; Bacteroidetes; Flavobacteriia; Flavobacteriales; Flavobacteriaceae; Flavobacterium

Package [MIGS: cultured bacteria/archaea\\_soil; version 5.0](#)

#### Attributes

<b>strain</b>	P4023
<b>collection date</b>	2011-02-02
<b>depth</b>	0
<b>elevation</b>	1630 m
<b>broad-scale environmental context</b>	polar environment [ENVO:01001703]
<b>local-scale environmental context</b>	polar biome [ENVO:01000339]
<b>environmental medium</b>	nest of bird [ENVO:00005805]
<b>geographic location</b>	<a href="#">Antarctica</a>
<b>isolation and growth condition</b>	R2A agar, 20 degrees C
<b>latitude and longitude</b>	<a href="#">63.77 S 57.78 W</a>
<b>number of replicons</b>	1
<b>reference for biomaterial</b>	CCM 9062T
<b>isolation source</b>	sample taken from an organic matter in an abandoned nest of a bird
<b>type-material</b>	type strain of Flavobacterium flabelliforme

Description Keywords: GSC:MlxS;MIGS:5.0

BioProject [PRJNA713758](#) Flavobacterium flabelliforme strain:P4023  
Retrieve [all samples](#) from this project

Submission Masaryk University, Stanislava Kralova; 2021-03-11



# Metadata v biologii

**Tvorba a uchování metadat – jedno pravidlo: „There is no good reason not to save metadata“**

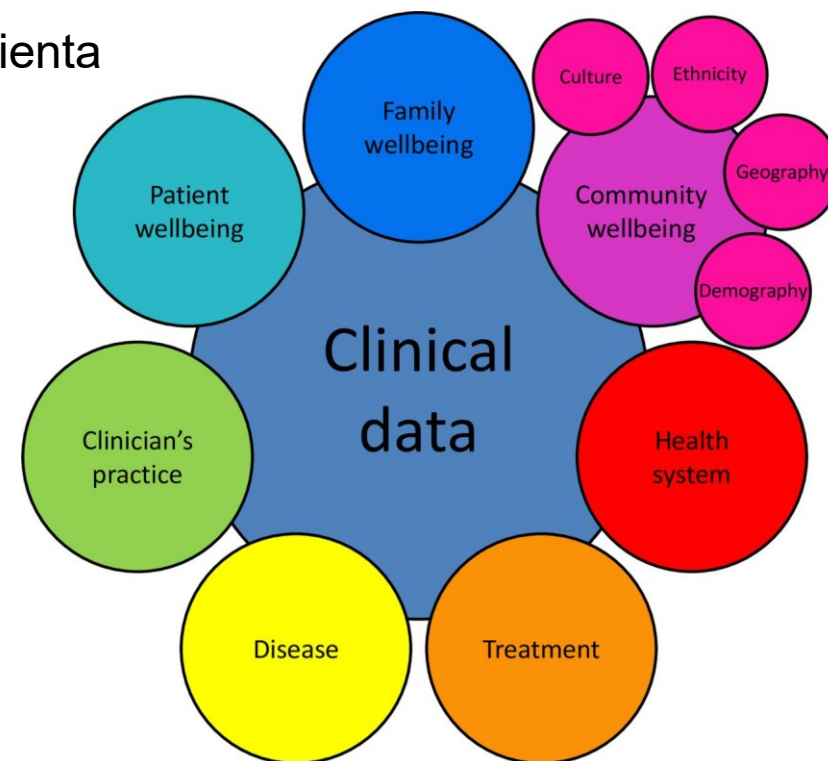
## EU projekty

- musíte uchovávat **vše**  
(*musíte mít odhad o velikosti a typu dat*)
- po dobu projektu + **5 let** po ukončení  
(*nutno přemyslet dopředu, kdo bude odpovědná osoba?*)
- aktuálně: metadata musíte nahrát **do veřejně dostupných** databází a spárovat s daty



# A co klinická data ve výzkumu?

- Specifické
- **Přísná pravidla** na sdílení informací
- **Pseudonymizace** – identifikace pacienta + dat možná dle klíče ( př. náhodné čísla jako ID)
- **Anonymizace (de-identifikace)** → oddělení dat od identifikace pacienta
- Neanonymizovaná data nesmí opustit nemocnici

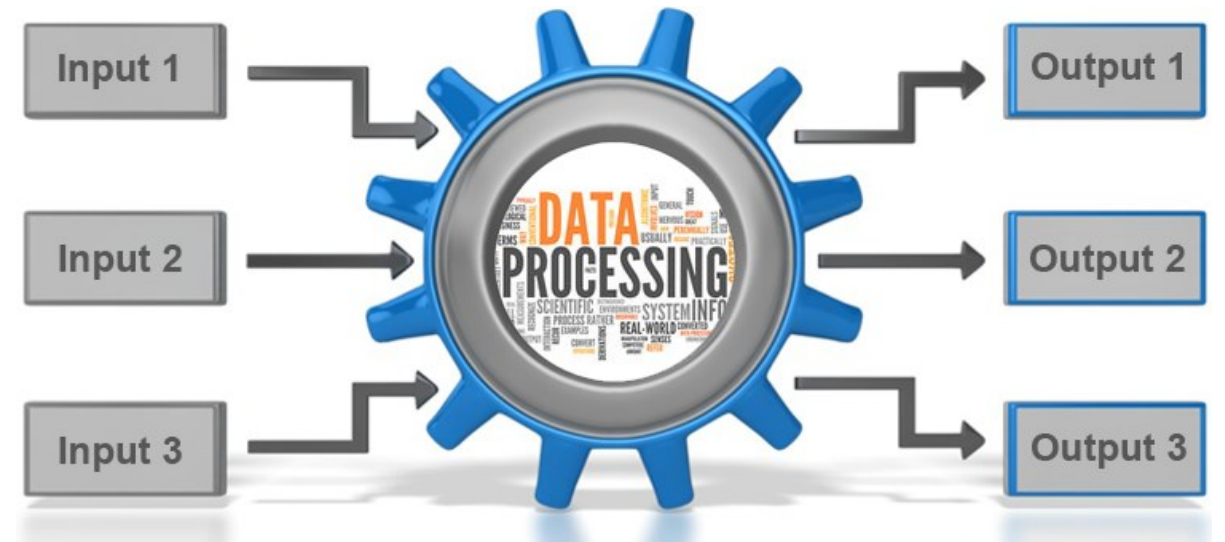




# Data processing



- Vložená nasbíraná data „input“
- Příprava dat (filtrace, kontrola kvality)
- Zpracování (z raw, tedy tzv. surových dat)
- Interpretace „output“



raw data  
často nečitelná člověkem

jednodušší  
zobrazení



# Krok 1: Filtrace dat

RAW DATA → ZPRACOVANÁ DATA

- Kvalitativní vyhodnocení dat
- **Normalizace** získaných dat pro možnost porovnání
- **Odstranění** chybných/nevyhovujících dat
- **Organizace**
- **Anotace**





# Příklady – filtrace dat

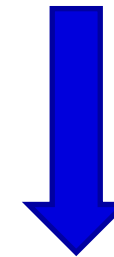
## Genomická data

### – sekvenace genomu

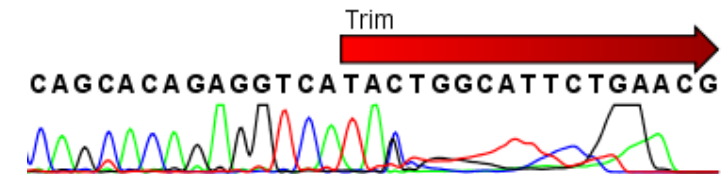
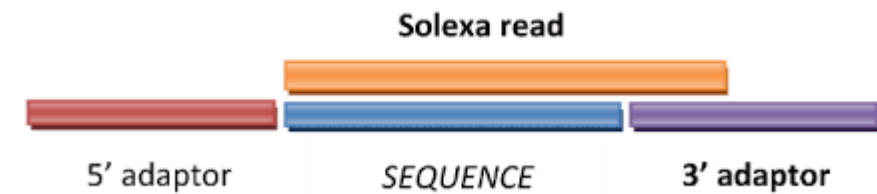
- výsledek tisíce čtení (readů) – nejsou stejné
- jiná délka
- jiná kvalita bází
- přítomnost adaptorů
- přítomnost kontaminací
- nekvalitní koncové sekvence (jiné báze než v skutečnosti)



Kontrola kvality  
→ filtrace



Výsledný soubor



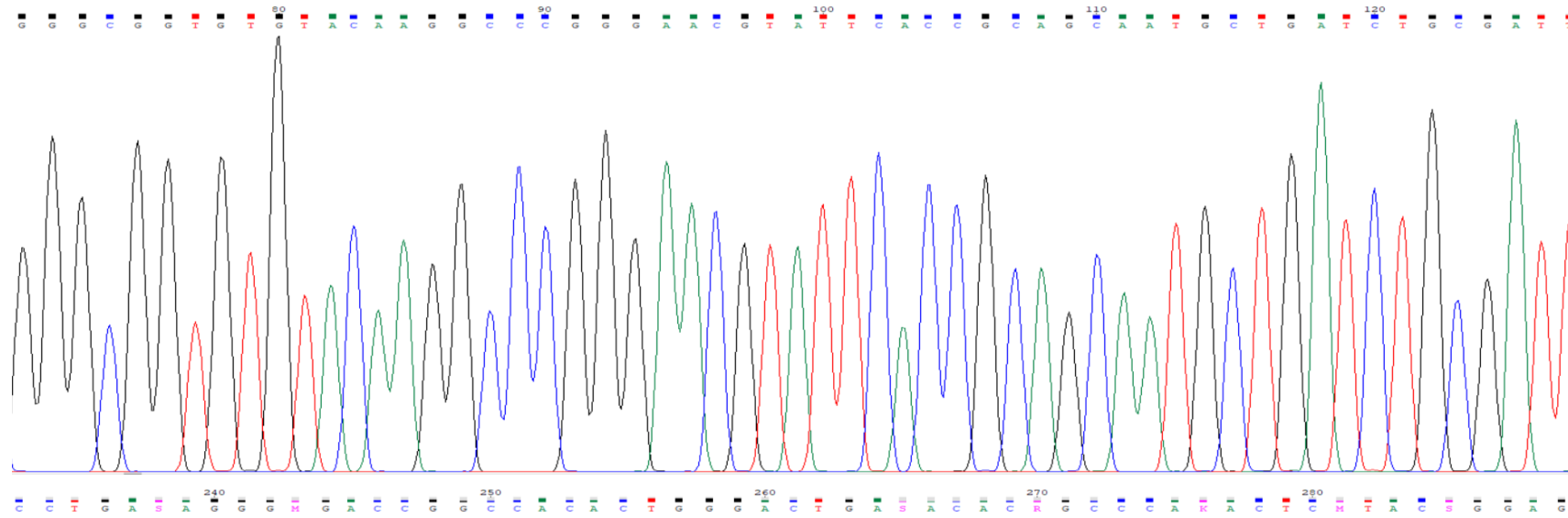
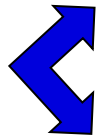
Phred skóre kvality –  
kvalita jednotlivých bází

Minimální délka readů

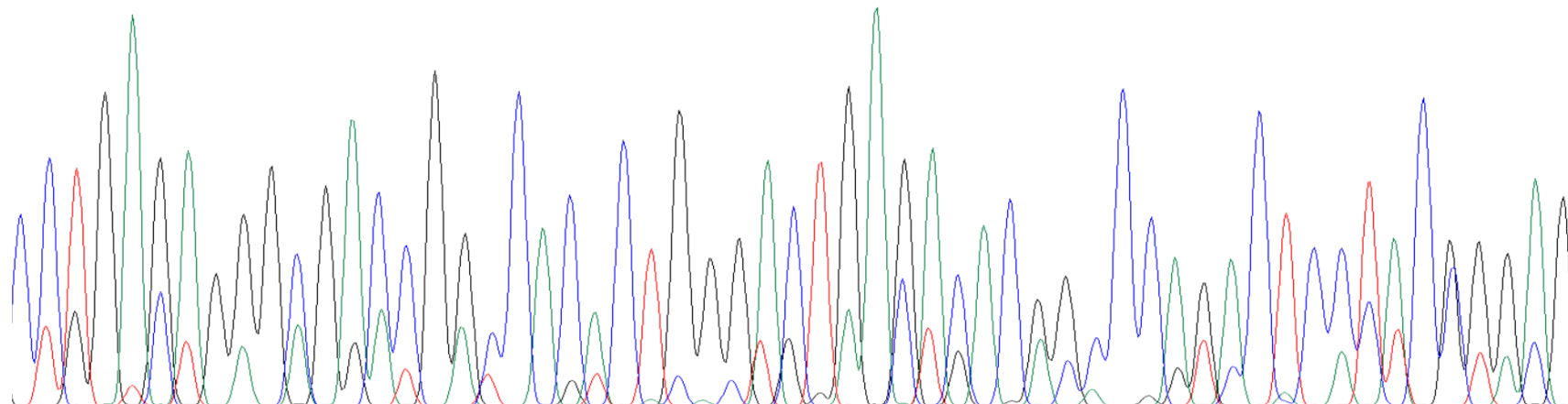


# Příklady – filtrace dat

čistá  
sekvence

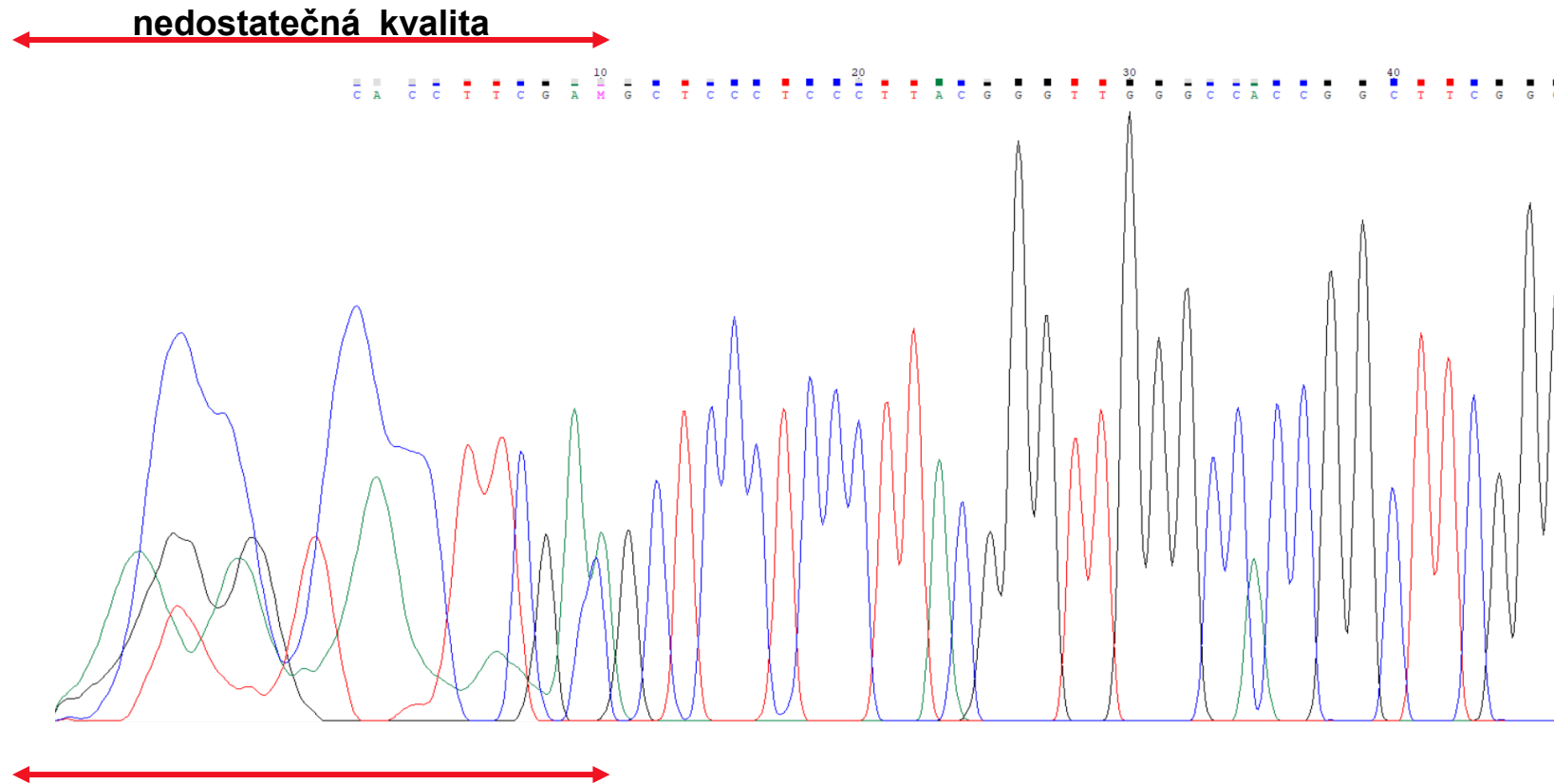


smíšená  
sekvence





# Příklady – filtrace dat

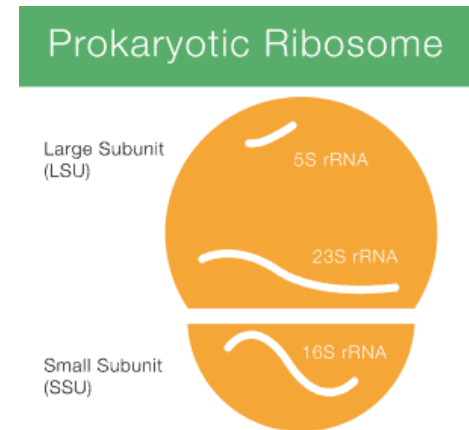




# Příklady – filtrace dat

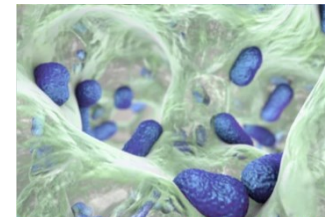
## Metagenom – lidský střevní mikrobiom

- PCR – výběr pouze jednoho genu (16S rRNA)
- sekvenace
- výsledek je nutné **filtrovat** (obdobně jako příklad 1)
- výsledek je nutné **normalizovat**



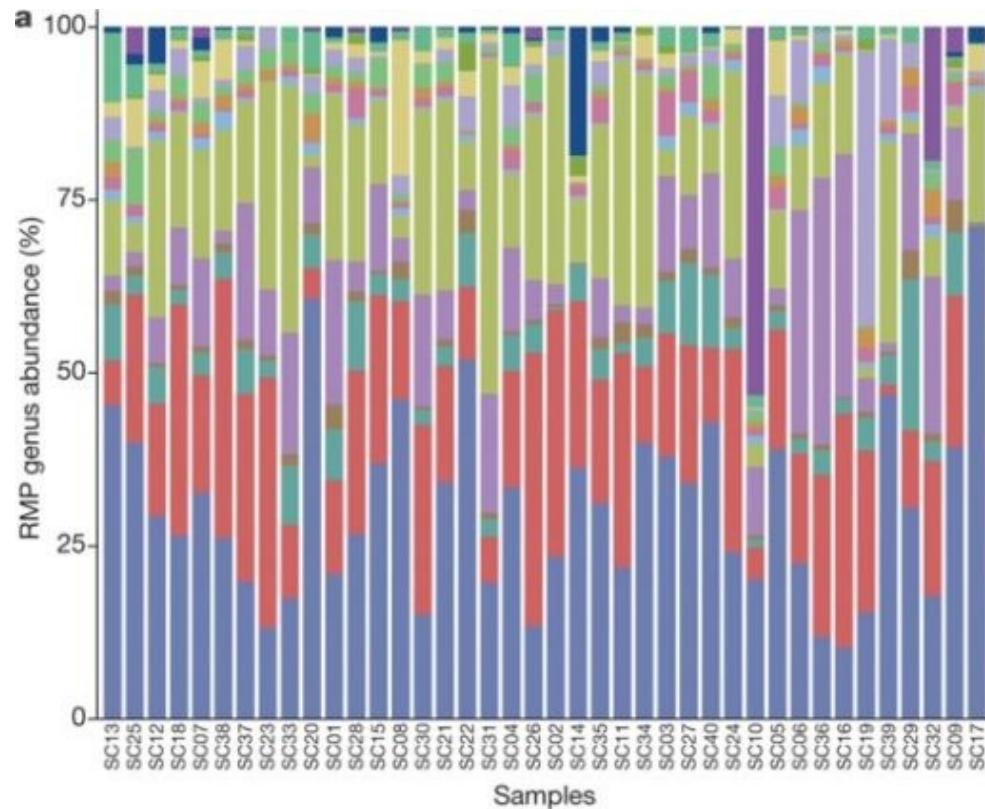
## příklad – výsledek tvrdí že máme v stolici:

- 100× *Bacteroides fragilis* a 100× *Prevotella copri* → 1:1
- normalizace genem 16S rRNA – *B.fragilis* 6 kopií X *P. copri* 4 kopie → 3:2

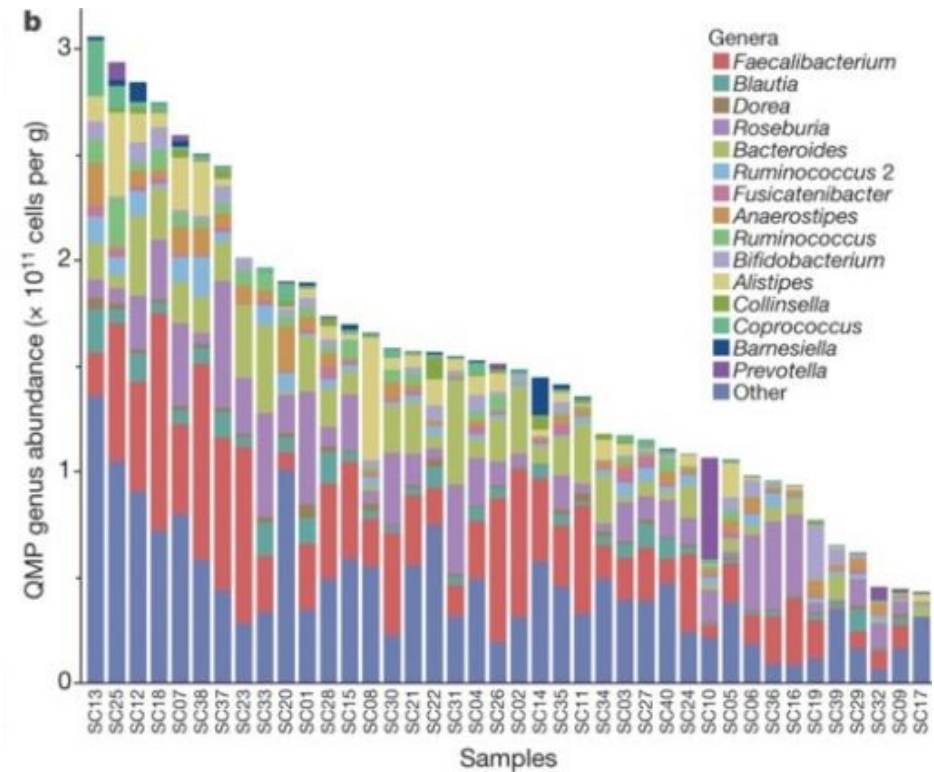




# Příklad normalizace sekvenčních dat



bez normalizace  
tzv „relativní početnosti“



po normalizaci vzhledem k velikosti vzorku  
tzv „absolutní početnosti“



# Kontrola kvality dat

## – Stanovení postupu kontroly kvality dat

*(kvalita fotek, kvalita sekvencí, úplnost metadat, kompatibilita dat, odstranění chybějících hodnot)*

## – Parametry kontroly vždy dle daných dat

- **sekvenace:** hloubka čtení, délka čtení, překrytí, kontaminace...
  - **experimentální data:** pozitivní a negativní kontroly, minimální počet opakování, reprodukovatelnost, statistické vyhodnocení...
- kvalita kontroly dat je nutná prakticky v každém kroku – sběr, generace dat, zpracování dat
- možná automatizace

# Problémy se zpracováváním dat

- Zejména BIOINFORMATIKA → tj. sekvenační data
- neexistují robustní univerzální pipeline
- každý si vytváří vlastní systém:
  - **reprodukovatelnost = těžká až nemožná**
  - **jedna data = různé výsledky!!!**





# Analýza dat



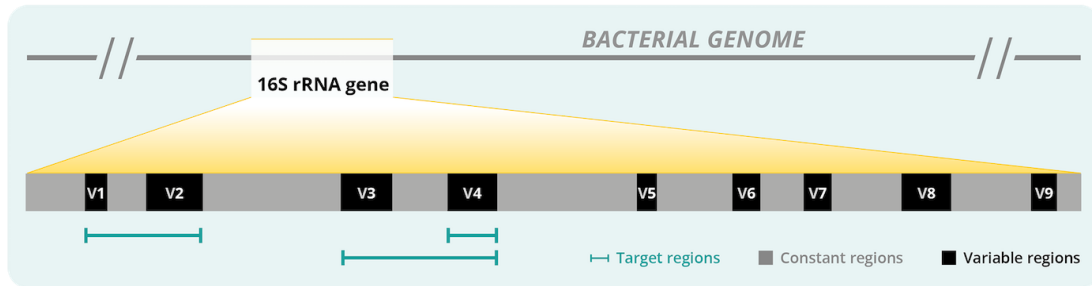
- biochemické a fyziologické data + způsoby vyhodnocení
- genomická data – bioinformatika (různé cíle = různé zpracování + příklady)
- metagenomická data (dtto)
- analýza dat od pacientu a asociace s výzkumnými daty

## Příklady:

- popisy nových bakteriálních druhů, fylogeneze, genomika
- sledování epidemiologie klinických kmenů
- hledání nových antibiotik

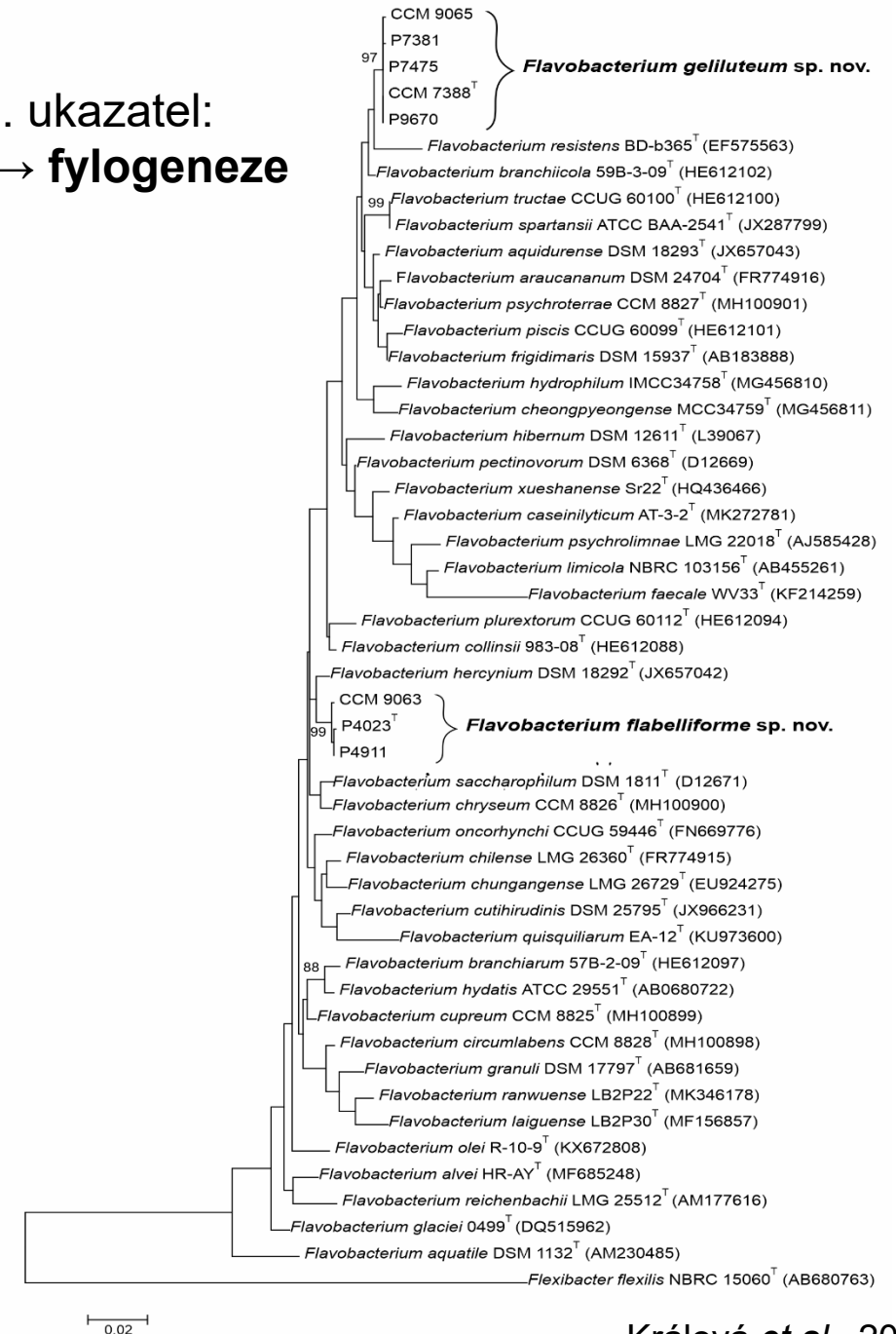
# Příklad analýzy 1 (nové bakteriální druhy – fylogeneze)

– Vstup: filtrovaná sekvenační data  
jednotlivé geny



1. ukazatel:  
→ **Cut-off value:** 98.65% podobnosti genu 16S rRNA

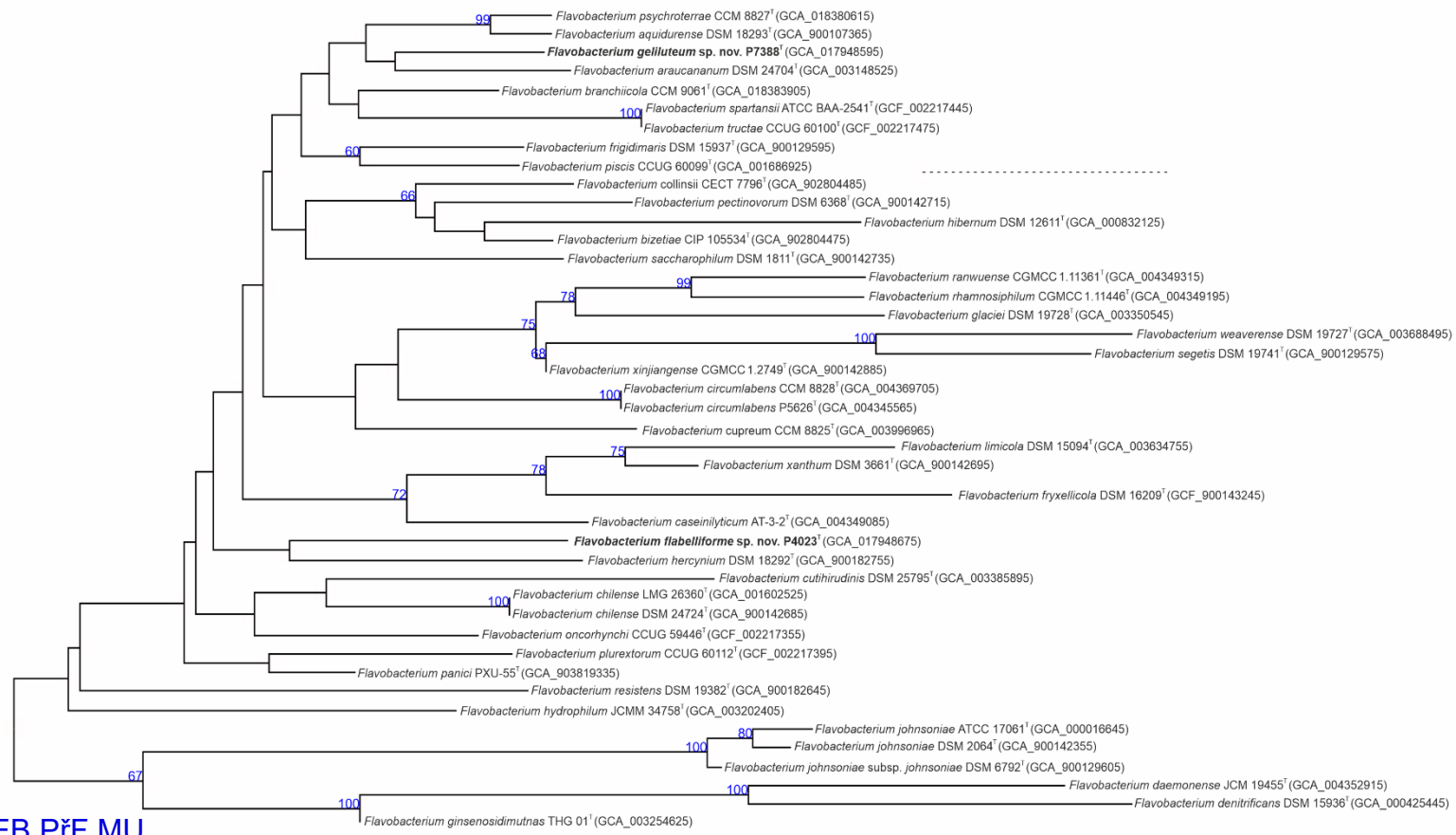
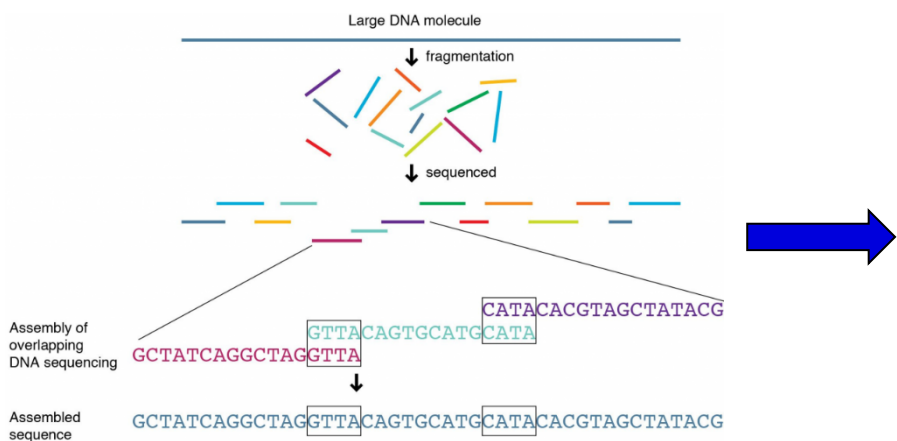
2. ukazatel:  
→ **fylogeneze**





# Příklad analýzy 1 (genomika)

– Vstup: filtrovaná sekvenační data, po kvalitě kontroly



1. ukazatel:

→ **Cut-off value:** 95% podobnosti genomu

2. ukazatel:

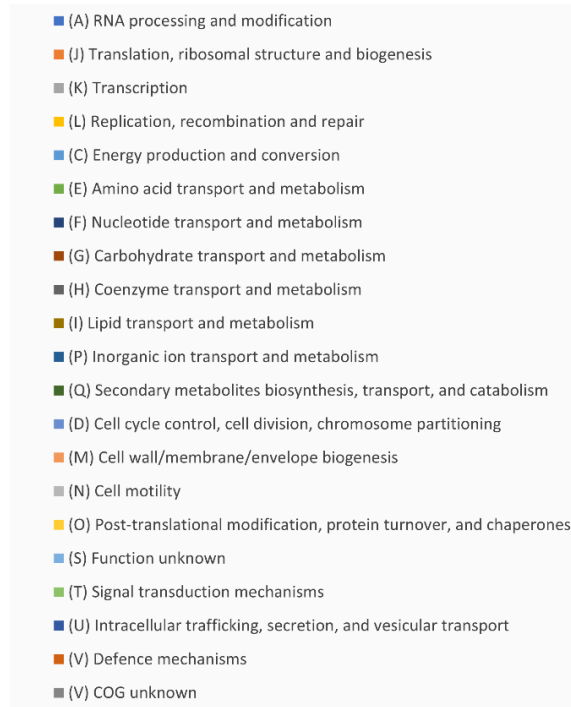
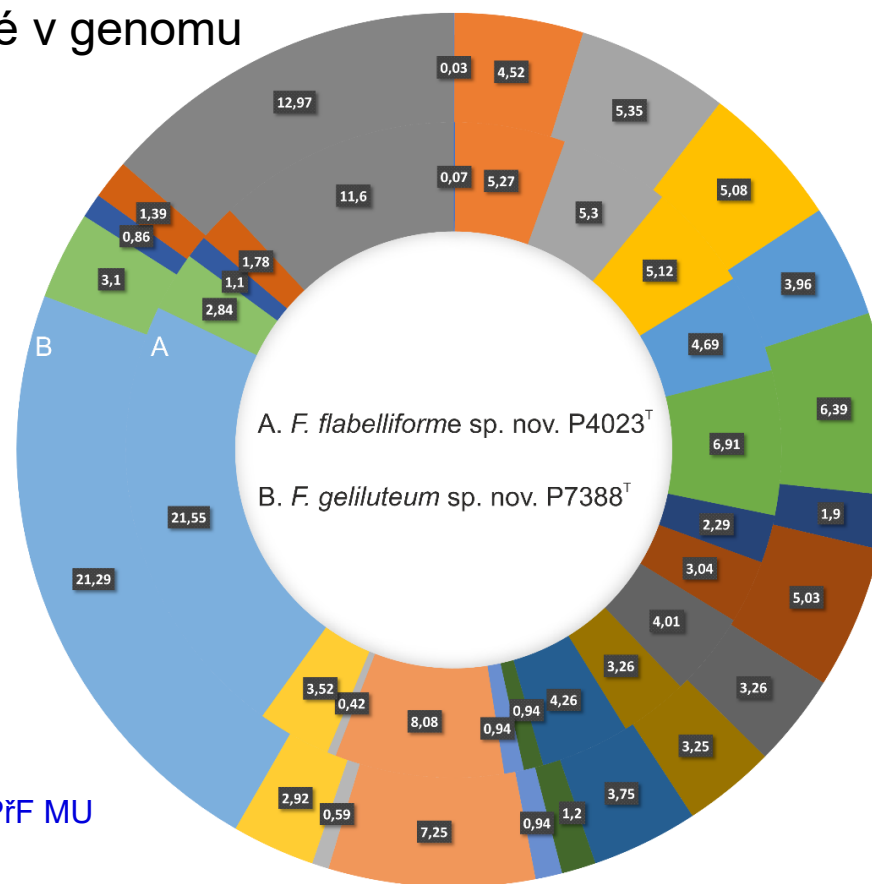
→ **fylogeneze**



# Příklad analýzy 1 (genomika)

– Vstup: filtrovaná sekvenační data, po kvalitě kontroly

- nejenom, jestli je bakterie „nová“
- můžeme přidat vlastnosti kódované v genomu





# Příklad analýzy 2 (epidemiologie)

## Stanovení příbuznosti izolátů

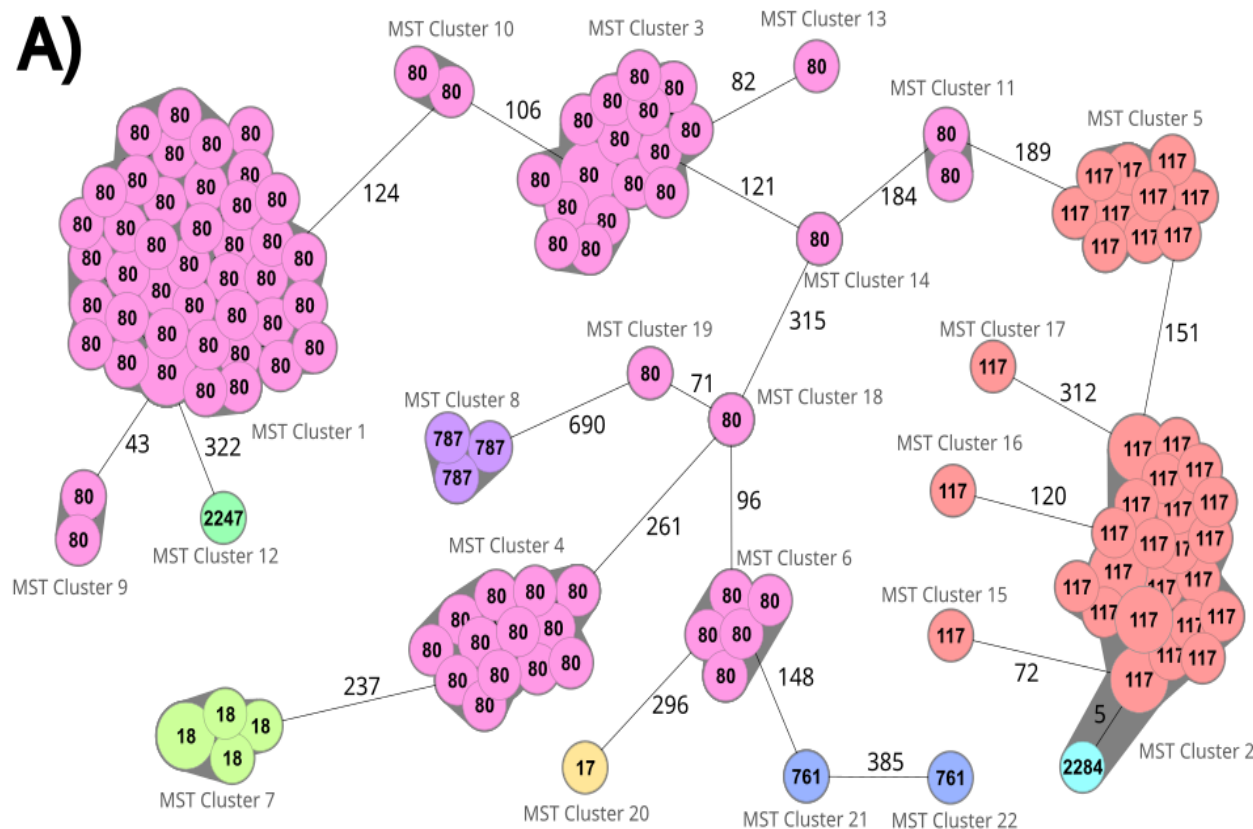
### VRE na základě WGS

#### – Vstup

- Složené genomy
- Metadata

#### – Výstupy

- Rozdělení do clusterů na základě počtu rozdílů v alelách jednotlivých genů
- Srovnání izolátů v rámci clusteru s dostupnými metadaty
- Na základě výsledků případně zahájení epidemiologického šetření a zavedení příslušných opatření



# Příklad analýzy 3 (biosyntetika)



- **Vstup: filtrovaná sekvenační data, po kvalitě kontroly**
- Vstup: fotografie
- Vstup: chromatografická data

# Příklad analýzy 3 (biosyntetika)



Select genome region:



## Identified secondary metabolite regions using strictness 'relaxed'

c00005\_NZ\_JABJ.. (original name was: NZ\_JABJRC010000014.1)



Region	Type	From	To	Most similar known cluster	Similarity
Region 5.1	siderophore	53,440	65,008		

c00011\_NZ\_JABJ.. (original name was: NZ\_JABJRC010000001.1)



Region	Type	From	To	Most similar known cluster	Similarity
Region 11.1	transAT-PKS, T1PKS, NRPS	586,329	666,426	9-methylstreptimidone	Polyketide:Modular type I 9%
Region 11.2	T3PKS	671,842	712,894	alkylresorcinol	Polyketide 100%
Region 11.3	melanin	1,268,489	1,278,860		
Region 11.4	siderophore	1,349,488	1,364,055		
Region 11.5	terpene	1,366,866	1,389,043	geosmin	Terpene 100%
Region 11.6	RRE-containing	1,499,997	1,521,088	CDA1b / CDA2a / CDA2b / CDA3a / CDA3b / CDA4a / CDA4b	NRP:Ca+-dependent lipopeptide 17%

c00015\_NZ\_JABJ.. (original name was: NZ\_JABJRC010000002.1)



Region	Type	From	To	Most similar known cluster	Similarity
Region 15.1	NRPS-like, NRPS	380,365	485,540	thiocoraline, NRP:Cyclic depsipeptide	34%
Region 15.2	lanthipeptide-class-i	496,456	520,654		
Region 15.3	redox-cofactor	895,339	917,379	lankacidin C, NRP + Polyketide	20%

c00016\_NZ\_JABJ.. (original name was: NZ\_JABJRC010000003.1)



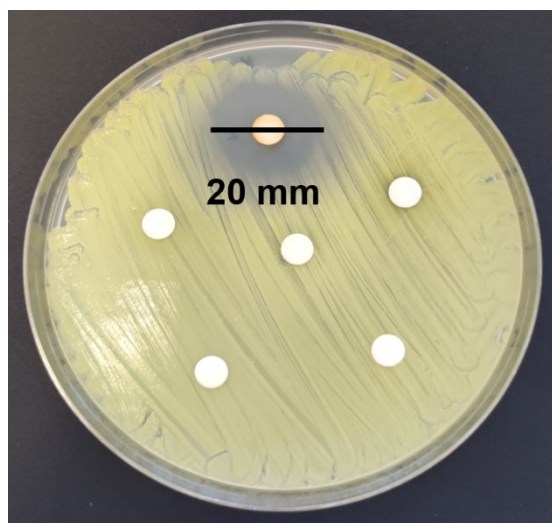
Region	Type	From	To	Most similar known cluster	Similarity
Region 16.1	lanthipeptide-class-iii	3,920	26,505	catenulipeptin, RiPP:Lanthipeptide	60%
Region 16.2	lanthipeptide-class-ii	31,433	54,363		



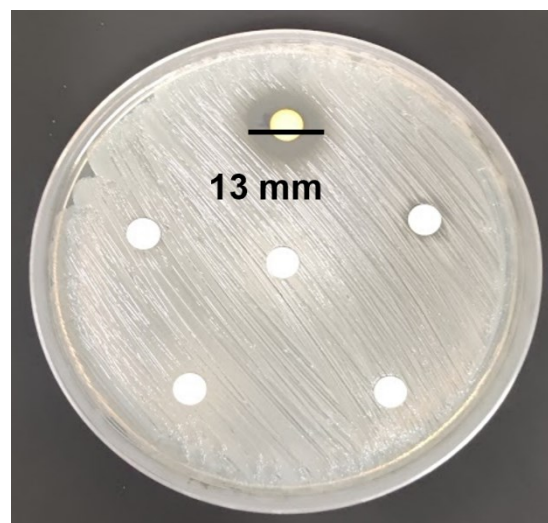


# Příklad analýzy 3 (biosyntetika)

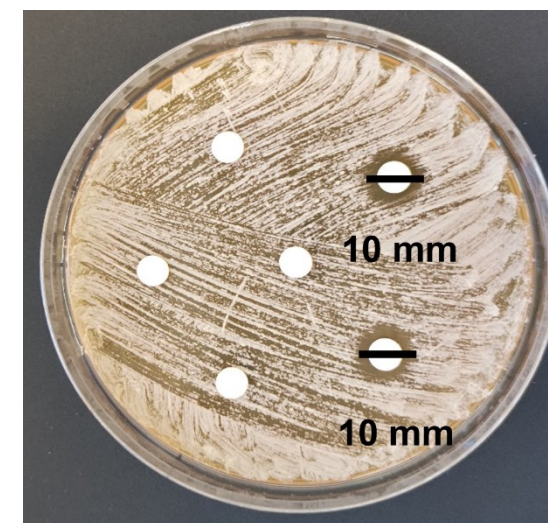
- Vstup: filtrovaná sekvenační data, po kontrole kvality
- **Vstup: fotografie**
- Vstup: chromatografická data



P12377, inhibition of *Micrococcus luteus*



P12377, inhibition of *Staphylococcus aureus*

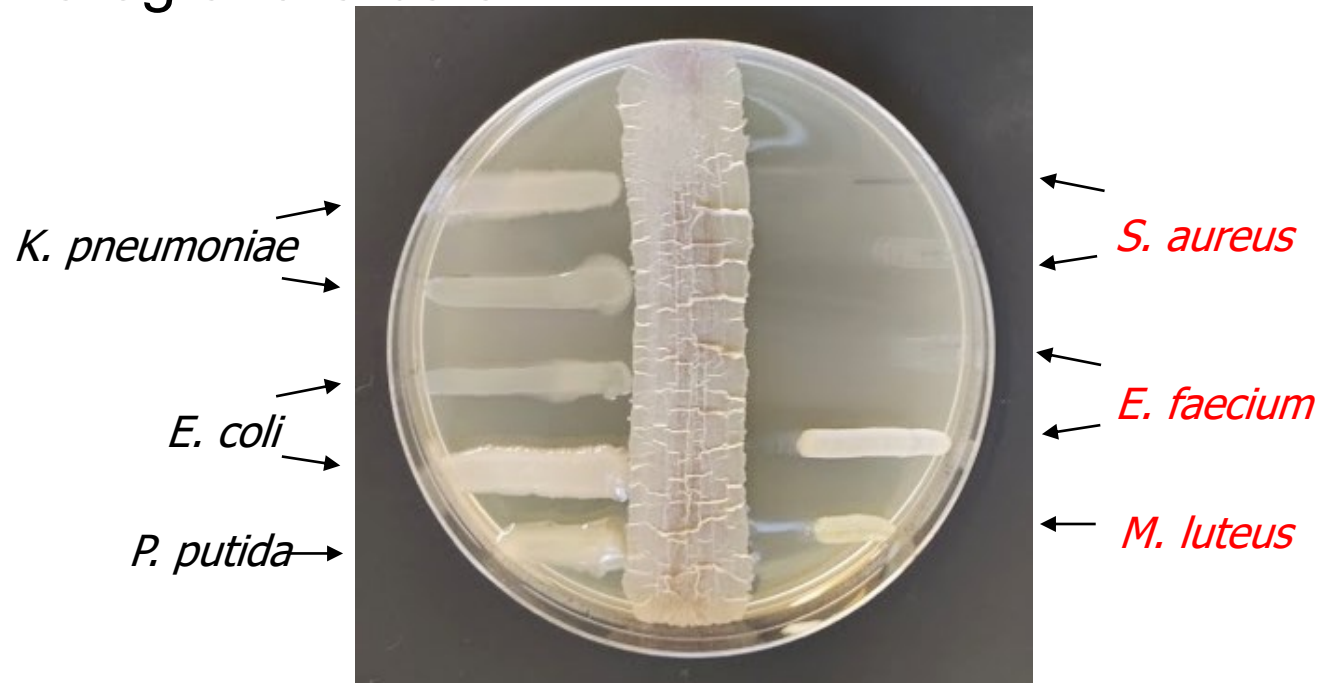


UV100, inhibition of *Saccharomyces cerevisiae*

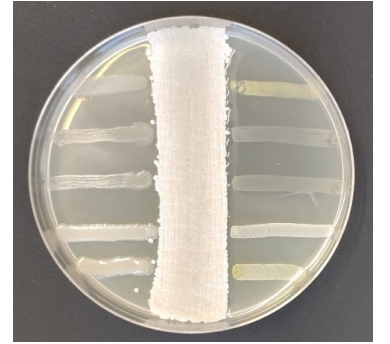


# Příklad analýzy 3 (biosyntetika)

- Vstup: filtrovaná sekvenační data, po kvalitě kontroly
- **Vstup: fotografie**
- Vstup: chromatografická data

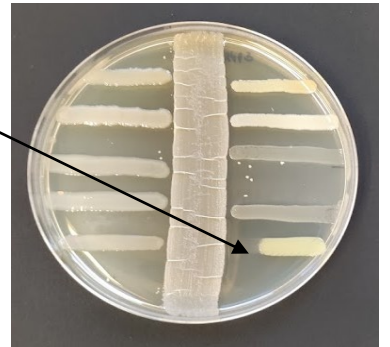


*Streptomyces* sp. UV100

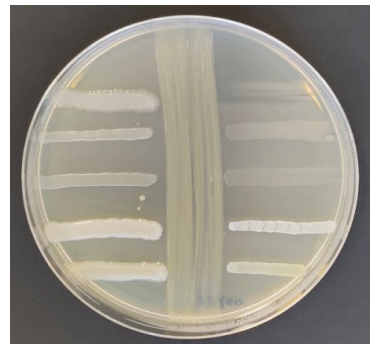


*Streptomyces* sp. P12413

*M. luteus*



*Arthrobacter* sp. P12200



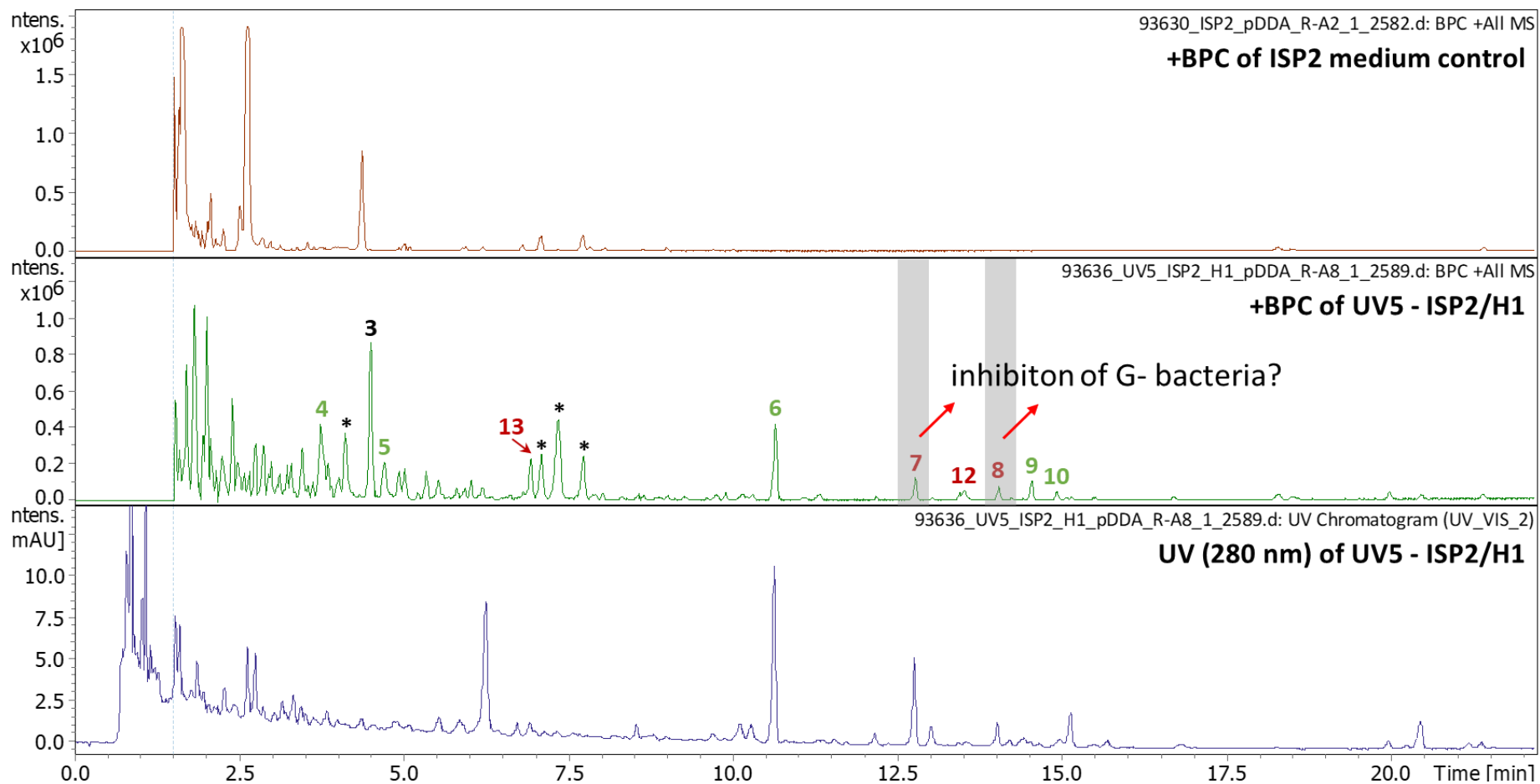
# Příklad analýzy 3 (biosyntetika)

- Vstup: filtrovaná sekvenační data, po kvalitě kontroly
- Vstup: fotografie
- **Vstup: chromatografická data**





# Příklad analýzy 3 (biosyntetika)





# Archivace dat



- data v tištěné podobě – doba uchování, katalogizace
- data v elektronické podobě
- lokální úložiště
- cloudové úložiště
- uzavřené databáze
- open access databáze
  
- **ovlivnění grantovými agenturami – podmínky + příklady**

# Archivace dat – příklad: srovnání univerzit



## Česká sbírka mikroorganismů, MU

- katalogizace v tištěné podobě
- katalogizace v digitální podobě, software MINE
- **výhody:**
  - dva systémy (záloha v tištěné podobě)
- **nevýhody:**
  - digitální systém zastaralý, neexistuje IT podpora
  - neexistuje záloha přístupná vzdáleně
  - neexistuje záloha přístupná všem

## CMESS, UniVie

- katalogizace/návody primárně digitální
- papírové formy pouze výjimečně (návodů k strojům)
- **výhody:**
  - univerzálně dostupném všem
  - pokud aktualizace, okamžitě u všech
  - vzdálený přístup samozřejmostí, pro všechny
  - denní zálohování na 3 serverech
- **nevýhody:**
  - nutný tisk pokud chcete vlastní kopii



# Archivace dat – online databáze a repozitáře

Ukládání dat online má v biologii pravidla:

- **NUTNOST** před každou publikací
- minimální doba – 5 let po publikaci
- perzistentní identifikátor (DataCite DOI)
- umožněno review ukládaných datasetů!
- ideálně specifické repozitáře dle povahy dat nebo dle standardu v daném oboru
- open access
  - (výjimka: klinická data mohou vyžadovat „*Data Usage Agreements*“)



# Archivace dat – příklady

## Data – nukleové kyseliny

Data types	Repository options	Data and metadata standards
Raw sequencing data (reads or traces) Genome assemblies Annotated sequences Sample metadata	INSDC repositories <a href="#">Genome Sequence Archive (GSA)</a>	Browse data and metadata standards endorsed by the Genome Standards Consortium
Genetic variation data	<a href="#">dbSNP</a> (human variations less than 50bp) <a href="#">dbVar</a> (human variations greater than 50bp) <a href="#">ClinVar</a> (human genotype & phenotype) <a href="#">European Variation Archive (EVA)</a> (all species) <a href="#">Genome Sequence Archive for Human (GSA-Human)</a>	

## Data – proteinové sekvence

<a href="#">UniProtKB</a>	<a href="#">view FAIRsharing entry</a>
---------------------------	--

## Data – funkční genomika

<a href="#">ArrayExpress</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Gene Expression Omnibus (GEO)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">GenomeRNAi</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">dbGAP</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">The European Genome-phenome Archive (EGA)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Database of Interacting Proteins (DIP)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">IntAct</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Japanese Genotype-phenotype Archive (JGA)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">NCBI PubChem BioAssay</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Genomic Expression Archive (GEA)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">GWAS Catalog</a>	<a href="#">view FAIRsharing entry</a>

## Data – obrázky

<a href="#">Image Data Resource</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">The Cancer Imaging Archive</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">SICAS Medical Image Repository</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Coherent X-ray Imaging Data Bank (CXIDB)</a>	<a href="#">view FAIRsharing entry</a>
<a href="#">Cell Image Library</a>	<a href="#">view FAIRsharing entry</a>

nature/repositories



# Archivace vzorků – mikrobiologie



## TAXONOMIE: živé kultury – reprezentující

nové druhy → musí být dostupné

- nejde ale o open access!
- uchovávaní zabezpečují sbírky mikroorganismů → poskytování je placené

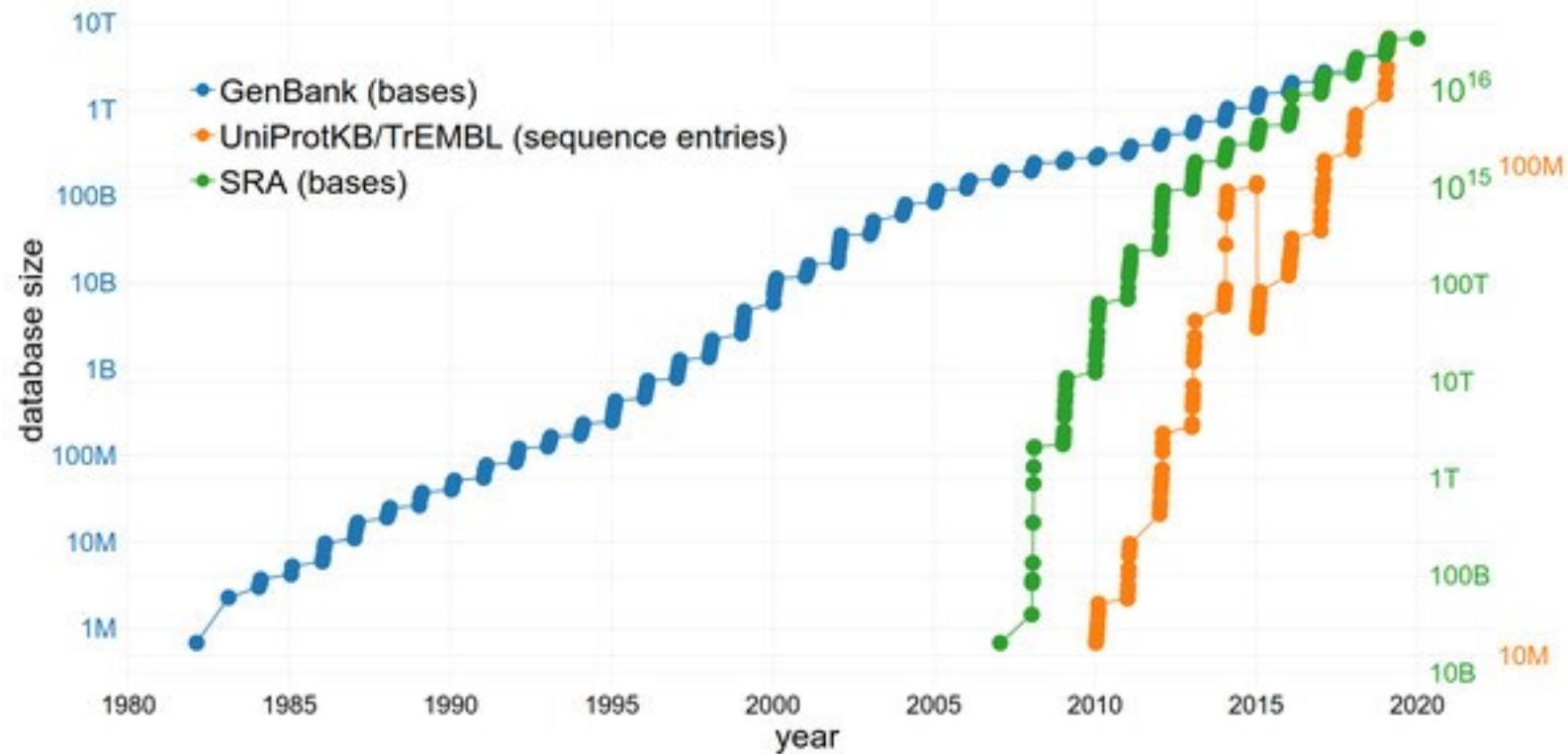
## BIODIVERZITA/EKOLOGIE: živé kultury

- nejde o nutnost, ale často jsou uložené volně dostupné soubory izolátů
- nejde/jde o open access
- uchovávaní zabezpečují sbírky mikroorganismů / university → placené/neplacené





# Data sharing



Sielemann *et al.*, 2020



# Data sharing – význam

## V biologii nutnost:

- poskytnout **všechna data**
- nebo poskytnou „**minimální dataset**“ → replikace výsledků (všechny data z článku, metadata, metody)  
„*Data Availability Statement*“ při submitování článku
- bez restrikce přístupu – open access
- pokud restrikce → etické/legální → jiný přístup **nutný** (agreement, official request, official permit)

# Data sharing



## V biologii nutnost:

- validace výsledků
- replikace, re-analýza, re-interpretace
- nové analýzy
- reprodukovatelnost
- archivace dat → investice do výzkumu má dlouhodobé výsledky
- citace dat → viditelnost a rozeznatelnost autorů, tvůrců dat, kurátorů





# Data sharing – význam i v praxi!

## Sledování, reakce, predikce:

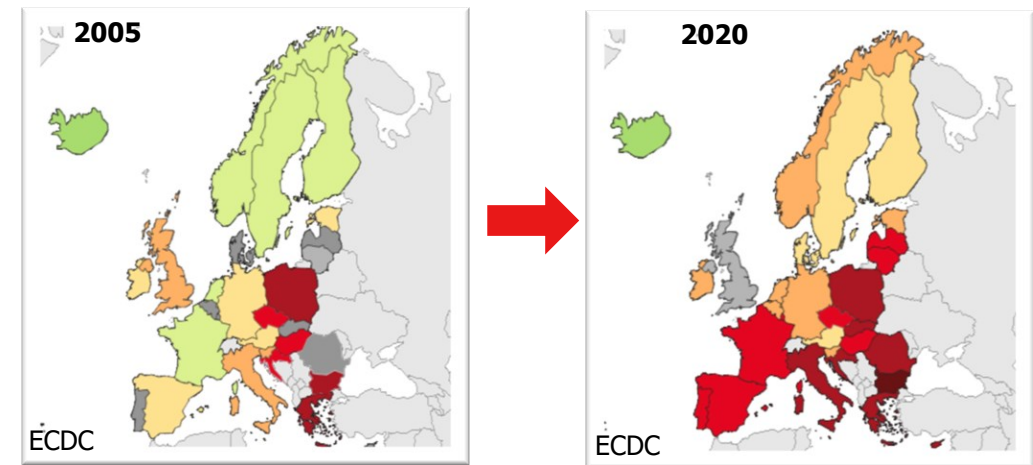
- infekční onemocnění (SARS-Cov-19, HAV v Brně, opičí neštovice...)
- šíření rezistentních patogenů (nemocnice, čističky vod...)
- šíření rezistence

## Zrychlení výzkumu:

- infekční onemocnění
- mikrobiom a související onemocnění
- Human Genome Project
  - rakovina
  - vzácná genetická onemocnění
  - farmakogenomika



Šíření rezistence k cefalosporinům v Evropě – *Klebsiella pneumoniae*







# Data reusability

- Použití **dříve nasbíraných/zpracovaných** dat za jiným účelem oproti původnímu plánu
- **Podmínka:**
  - vysoká kvalita
  - eticky v pořádku
  - dostatečně popsané – metadata
  - aktuální verze dat
  - uvedené podmínky opakovaného použití
  - **NUTNO OCITOVAT, UVÉST REFERENCI!!!**





# Data reusability – význam

- nezávislé studium dat různými skupinami/vědci
- šetření finančních zdrojů
- urychlení výzkumu
- násobné možnosti zpracování a vyhodnocení (jedny data = mnoho výsledků)
  
- **Nevýhody:**
  - neúplné/nesprávné metadata
  - nízká kvalita dat
  - chybná data



**Sequences: text strings describing sequential bases, including gaps in data and their length**  
(single gene/whole genome sequences, amino acid sequences)

```
>seq1
ATCGTTTAGCTAGACCTGATG
ATCCGATCGATTACGTG
>seq2
GACACGATCGTCAGAAATGCA
GTC
>seq3
ACGACAAATCATCTCC
>seq1
MYVRANQEFFK
>seq2
WTSMDCHLKV
>seq3
MGPKLHIGRQEFKLIHYWN
NNG
```

**Metadata: information about the acquisition, processing and presentation of data**  
(methodology, sample size and origin)

**Graphs: graphical representation indicating relationship**  
(pathway data, genetic maps, plasmid maps)



**Annotations (gene annotation, gene model)**

Primary data

**Algorithms and Software (code, bioinformatics pipelines)**

```
for ID in transcripts_per_genes[ gene ]:
    counter = 0
    for element in transcripts[ ID ]:
        if element[ 'type' ] == "CDS":
            counter += element[ 'end' ] -
            element[ 'start' ]
            CDS_len_per_transcript.append( { 'id': ID, 'len': counter } )
repr_trans = sorted( CDS_len_per_transcript, key=itemgetter( 'len' ) )
repr_transcripts.update( { repr_trans[ 'id' ]: transcripts[ repr_trans[ 'id' ] ] }
```

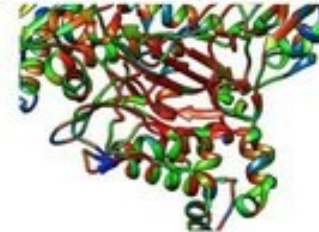
**Hypotheses, evidence and prose (data interpretation)**

**Plots and Images: natural, artificial and stylised imagery**  
(micrographs, radiographs, diagrams, microscopy of biological material, images of geographical regions)



**Measurement parameters**  
(enzyme affinity/speed, chromatography results, mass spectra, geo data/coordinates)

**Geometric information/Molecular structure data**  
(secondary and tertiary protein structures)



**Patterns**  
(sequence motifs, regulatory sequences, expression profiles)



**Primary databases: direct submission of experimentally-derived data from researchers; archival database**  
(DNA, RNA, protein, expression, disease, organism-specific databases)



Derived data

**Publications (text mining)**



**Meta-analyses**

**Derived databases: results of analysis, literature search and interpretation, curated database**  
(DNA, RNA, protein, expression, disease, organism-specific, phenotypic databases)



# Data reusability – příklad

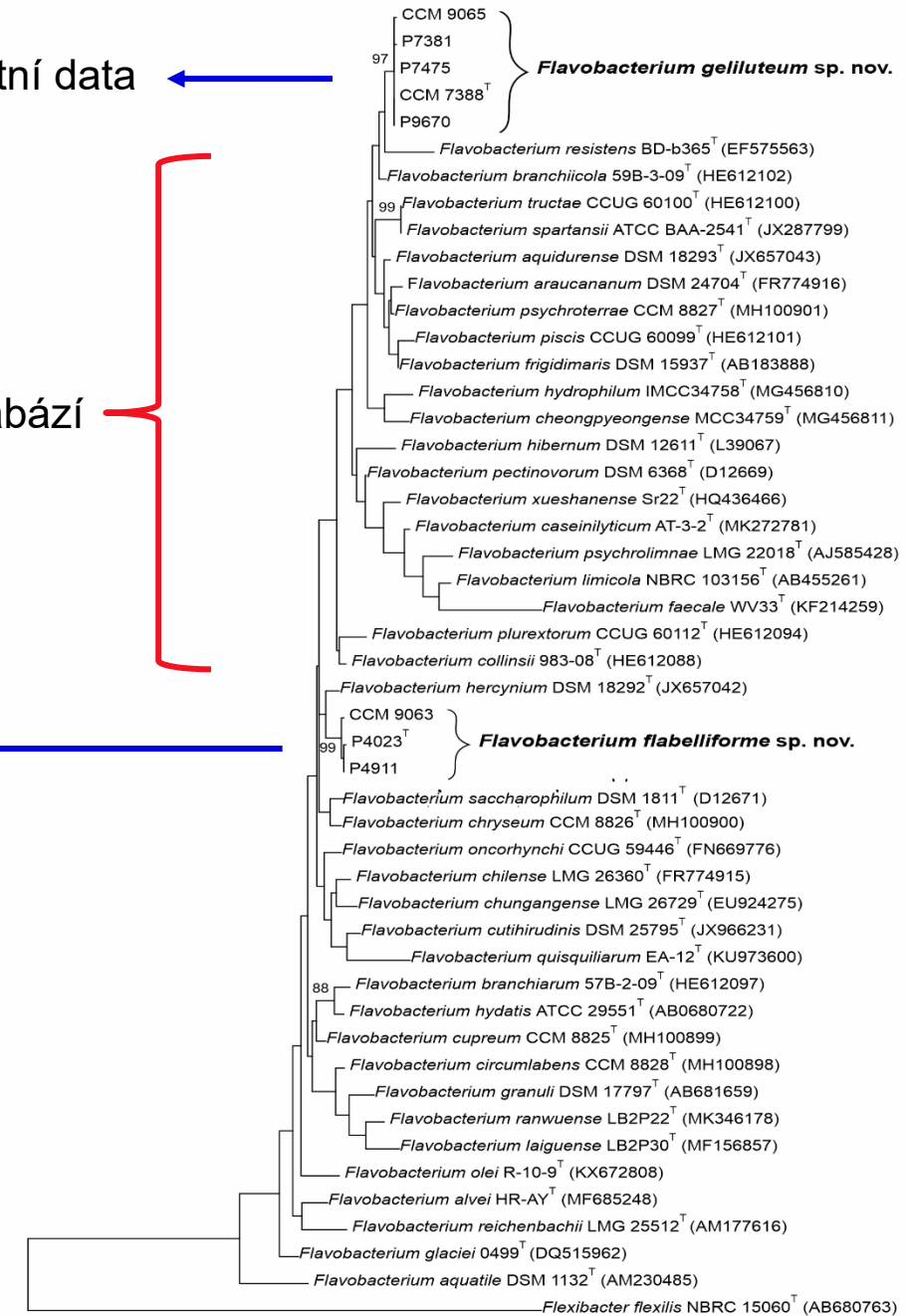
– každá fylogenetická analýza

data z repozitářů/databází

vlastní data

vlastní data

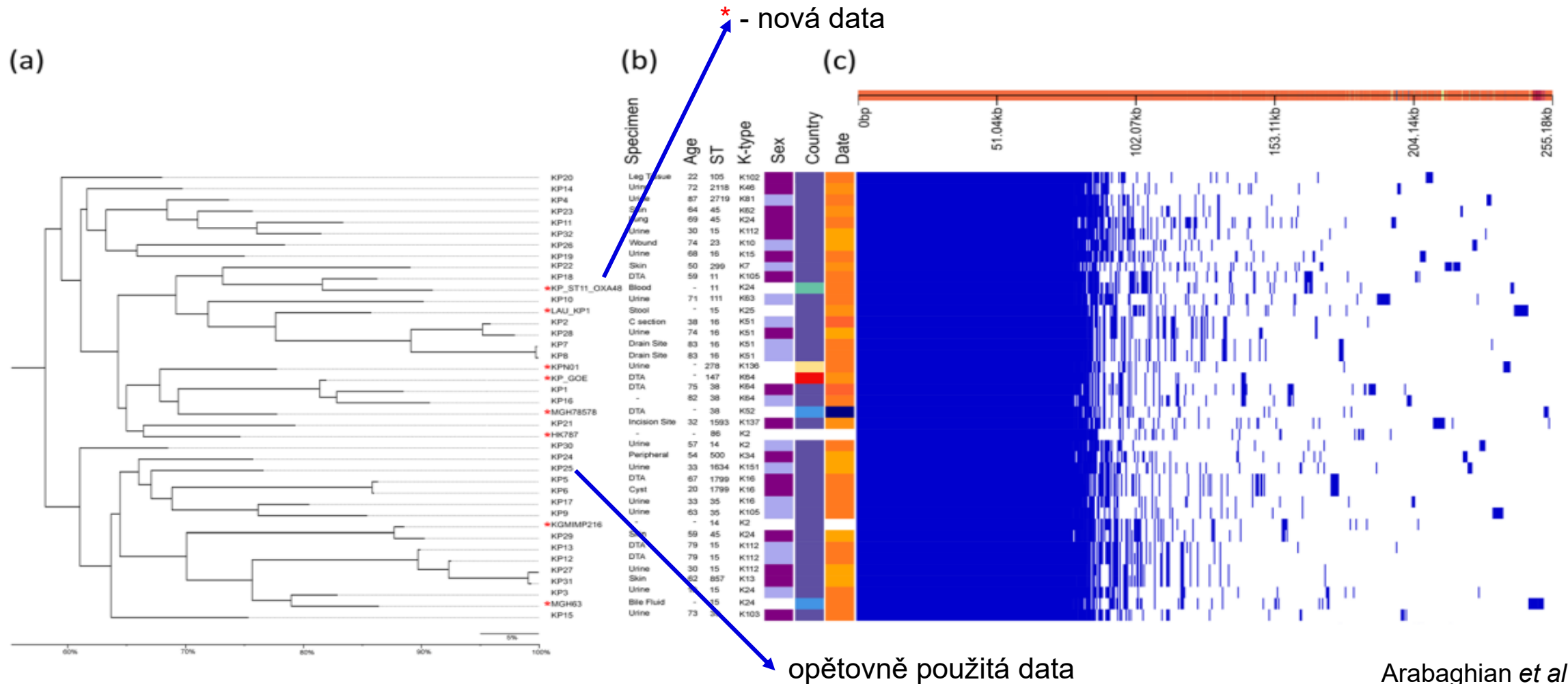
data z repozitářů/databází





# Data reusability – příklad

## Pangenomická analýza – *Klebsiella pneumoniae*



# Data reusability – příklad – klinický výzkum

## Tvorba epidemiologických schémat

- využití NCBI databází
- využití genómové subdatabáze
- možnost filtrování výsledků
- možnost stáhnutí celého souboru

The screenshot shows the NCBI Genome download interface. At the top, it displays the U.S. Department of Health and Human Services logo and the NIH National Library of Medicine logo. A search bar is present with the text "Search NCBI ...". Below the search bar, there are navigation links for "Datasets", "Taxonomy", "Genome" (which is highlighted), "Gene", "Command-line tools", and "Documentation". The main heading is "Genome" with a "BETA" badge. Below the heading, there is a description: "Download a genome data package including genome, transcript and protein sequence, annotation and a data report". A "Selected taxa" dropdown menu is set to "Enterococcus faecium". Below this, there are "Filters" and "SEARCH WITHIN RESULTS" sections. The "Filters" section includes checkboxes for "Reference genomes", "Annotated genomes" (with sub-options for "Annotated by NCBI RefSeq" and "Annotated by GenBank submitter"), and "Exclude atypical genomes". The "SEARCH WITHIN RESULTS" section has a search bar and a description: "Enter taxon name or modifier, assembly name or submitter". Below this, there are two sliders: "ASSEMBLY LEVEL" ranging from "contig" to "complete" and "YEAR RELEASED" ranging from "1980" to "2022".

# Data reusability – příklad – klinický výzkum

## Tvorba epidemiologických schémat

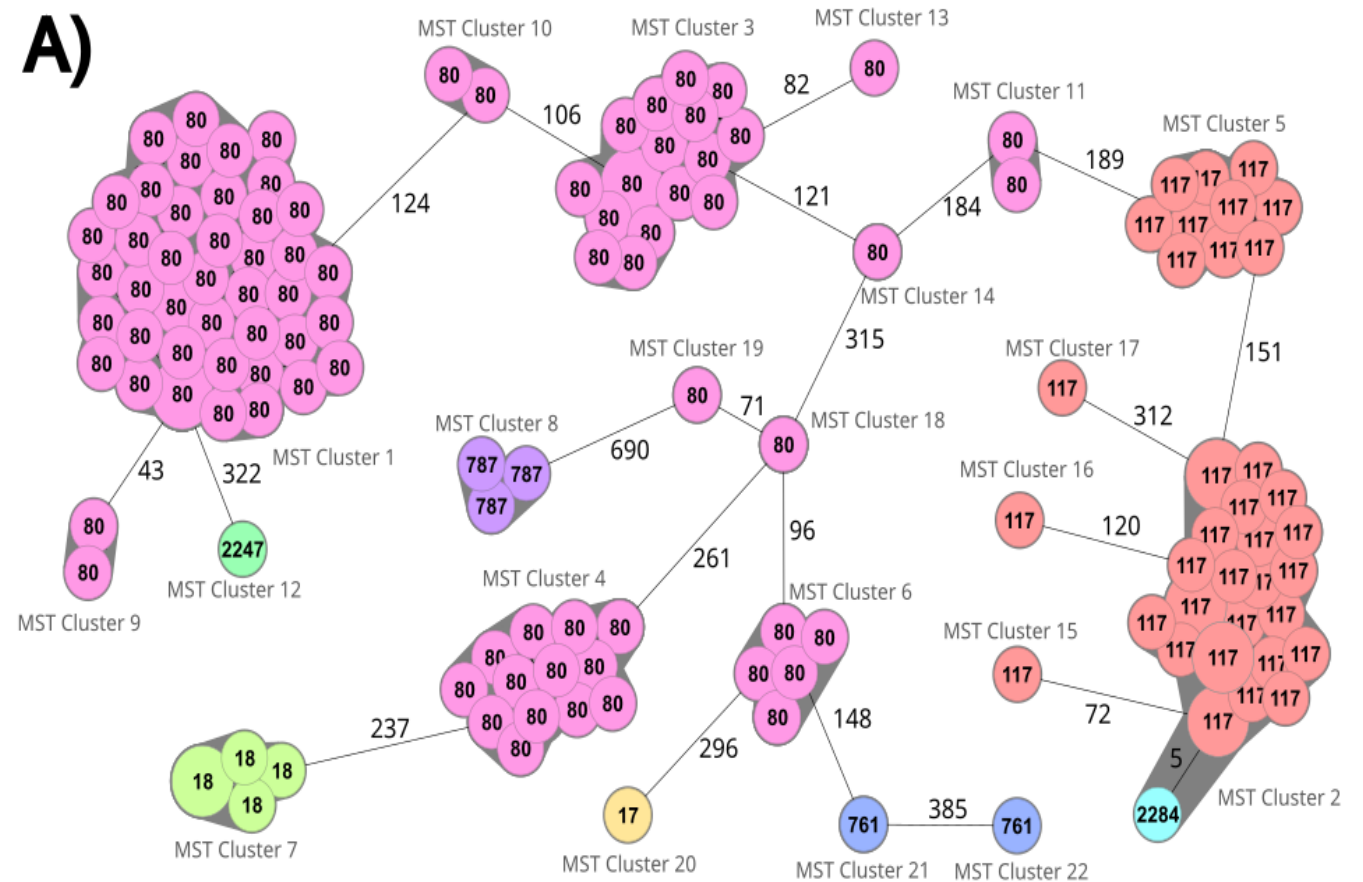
Download ▾ Select columns 17 390 genomes Rows per page 20 ▾ 1-20 of 17 390 < >

<input type="checkbox"/> Assembly	Scientific name	Modifier	Annotation	Size (Mb)	Level	Year	WGS acce	Action
<input type="checkbox"/> ASM973400v2 <span>reference</span> RefSeq: GCF_009734005.1 GenBank: GCA_009734005.2	Enterococcus faecium	SRR24 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.919	Complete	2020		⋮
<input type="checkbox"/> ASM76497v1 RefSeq: GCF_000764975.1 GenBank: GCA_000764975.1	Enterococcus faecium	UC7266 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.806	Contig	2014	JRJW01	⋮
<input type="checkbox"/> ASM76498v1 RefSeq: GCF_000764985.1 GenBank: GCA_000764985.1	Enterococcus faecium	UC7265 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.807	Contig	2014	JRHQ01	⋮
<input type="checkbox"/> ASM76734v1 RefSeq: GCF_000767345.1 GenBank: GCA_000767345.1	Enterococcus faecium	70-7-8 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.476	Scaffold	2014	JRUF01	⋮
<input type="checkbox"/> ASM76736v1 RefSeq: GCF_000767365.1 GenBank: GCA_000767365.1	Enterococcus faecium	70-36-8 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.717	Scaffold	2014	JRUG01	⋮
<input type="checkbox"/> ASM77250v1 RefSeq: GCF_000772505.1 GenBank: GCA_000772505.1	Enterococcus faecium	VRE3 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.820	Contig	2014	JSET01	⋮
<input type="checkbox"/> ASM77252v1 RefSeq: GCF_000772525.1 GenBank: GCA_000772525.1	Enterococcus faecium	ATCC 51559 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.954	Scaffold	2014	JSVT01	⋮
<input type="checkbox"/> ASM78705v1 RefSeq: GCF_000787055.1 GenBank: GCA_000787055.1	Enterococcus faecium	L-3 strain	<a href="#">NCBI RefS...</a>	2.642	Scaffold	2014	JRGX01	⋮
<input type="checkbox"/> ASM78706v1 RefSeq: GCF_000787065.1 GenBank: GCA_000787065.1	Enterococcus faecium	L-X strain	<a href="#">NCBI RefS...</a>	2.710	Contig	2014	JRGY01	⋮
<input type="checkbox"/> 70-61-7 RefSeq: GCF_000804385.1 GenBank: GCA_000804385.1	Enterococcus faecium	70-61-7 strain	<a href="#">NCBI RefS...</a> <a href="#">Submitter</a>	2.521	Scaffold	2014	JUEK01	⋮

# Data reusability – příklad – klinický výzkum

## Tvorba epidemiologických schémat

- příklad využití dat uložených za účelem publikace
- nový účel → epidemiologie



MUNI

# Shrnutí

6. přednáška CORE042

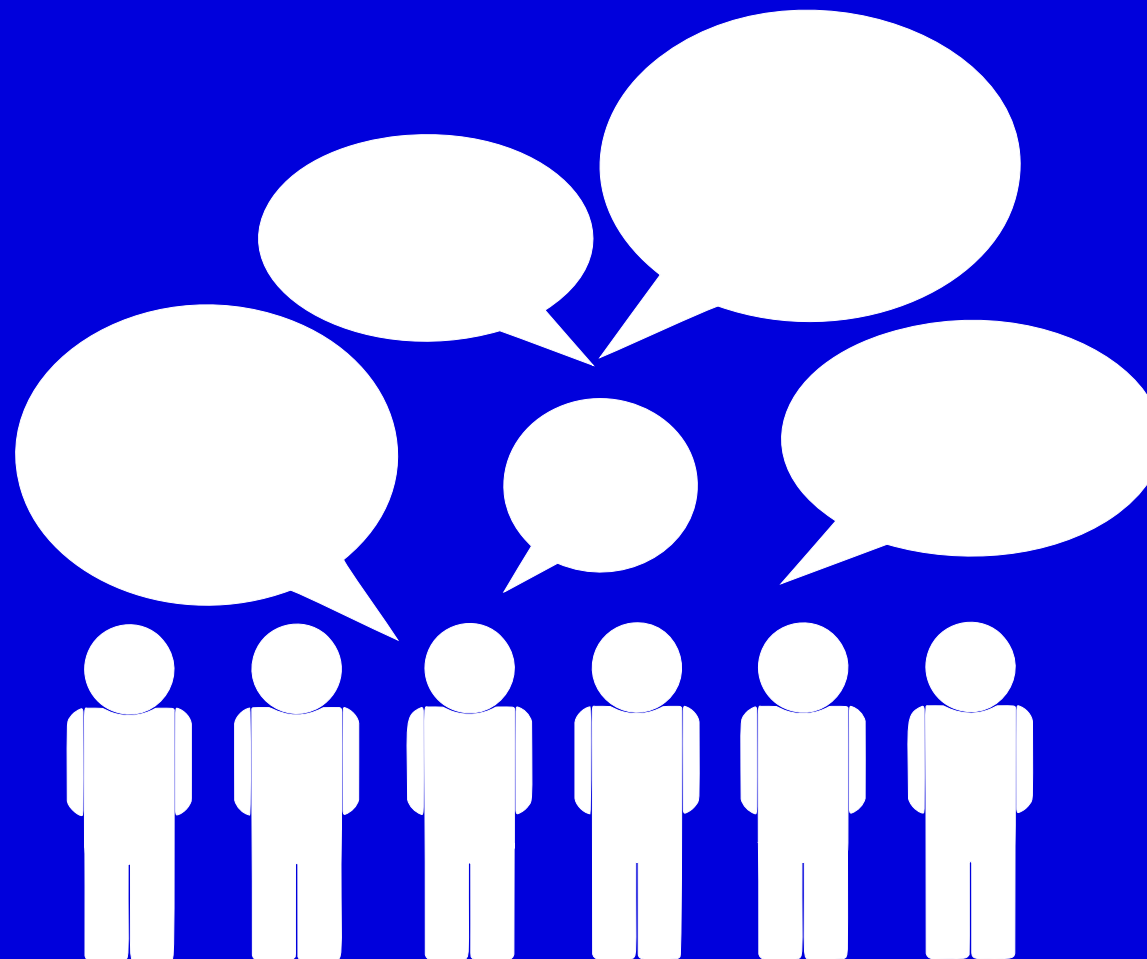


# Shrnutí

- Cyklus dat v mikrobiologii/biologii → využíváme všechny kroky
- Generujeme velké množství variabilních dat (fotografie, HPLC, sekvence, měření, popisky, ...)
- Generujeme enormní množství metadat (GPS souřadnice, typy materiálů, množství materiálů, časy/doby odběrů, klinická data, asociovaná data, ...)
- Vysoká náročnost na **plan** → **collect** → **process** → **analyse**
- Výrazný podíl **sharing + reusability**
- Výrazný tlak na FAIR data (zejména **preserve + share** → **to be reused!**)
  - velké množství databází
  - specifické databáze
  - nevýhoda – databáze bez review procesu a kontroly → chyby

# MUNI

## Diskuse



Zdroj: [Communicate\\_communication\\_conference\\_2028004](#) od [OpenClipart-Vectors](#) z [Pixabay](#)