# FAIR raziskovalni podatki v biologiji
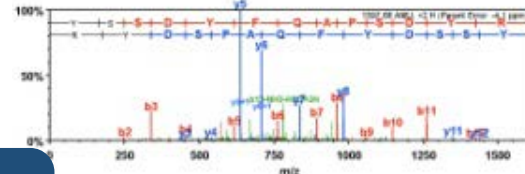
dr. Marko Petek

marko.petek@nib.si

12.05.2022

# Raziskave na Oddelku za biotehnologijo in sistemsko biologijo
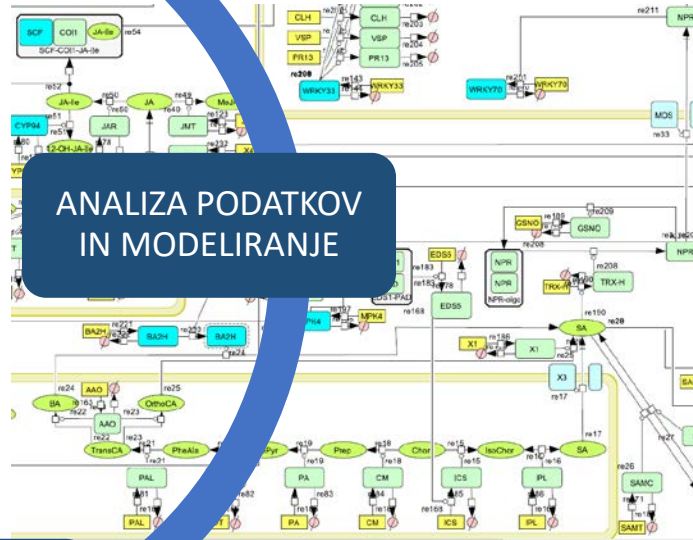
Delovni sklop: OMIKE
Vodja: prof. dr. Kristina Gruden



GENOMIKA, TRANSKRIPTOMIKA, PROTEOMIKA, METABOLOMIKA, …

POTRJEVANJE HIPOTEZ S FUNKCIJSKIMI ŠTUDIJAMI

ANALIZA PODATKOV IN MODELIRANJE

POSTAVLJANJE HIPOTEZ

# Kakšne raziskovane podatke generiramo?

- Surovi podatki
- Metapodatki
- Vmesni in končni rezultati obdelave podatkov in modeliranja
- Publikacije (članki, poročila, predstavitve, posterji etc)

# Načrt ravnanja s projektnimi podatki ("data management plan")

| | | | | | | |
|---|---|---|---|---|---|---|
| **Generated by** | CSIC, EI, TUDA | NIB, EI | TUDA, CSIC | TUDA | CSIC, EI | CSIC |
| **Format of generated data** | GenBank | fastq, hdf5 | DATA.MS, NMReDATA, Bruker data format, Shimadzu data format (qgd) | txt(ASCII), jpeg, eag | jpeg, txt(ASCII) | txt(ASCII), jpeg |
| **Raw data storage at** | CSIC, EI, TUDA, GeneBank | NIB, EI, SRA/ArrayExpress/GEO | TUDA, CSIC, MetaboLights, MolCheck | TUDA, FAIRDOMHub | CSIC, EI, FAIRDOMHub | CSIC, FAIRDOMHub |
| **Expected data size** | < 1 GB | < 1800 GB | < 30 GB | < 5 GB | < 10 GB | < 1 GB |
| **Analysed by** | CSIC, EI, TUDA | NIB, EI | TUDA, CSIC | TUDA | CSIC, EI | CSIC |
| **Analysed data storage at** | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | GB Elements Database, FAIRDOMHub |
| **Minimal information requirements** | MIRIAM | MIAME | CIMR | MINI | MIAPPE | / |
| **Standards, formats** | Genebank, SBOL data and visual | MAGE-ML | mzML, mzQuantML, nmrML | NWB | / | / |
| **Ontologies and vocabularies used** | SBOL | GO, KEGG, InterPro, MapMan | CHEBI, KEGG, MapMan | OEN | PO, TO, CO | SBOL |
| **SOPs stored at** | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub |
| **Scripts stored at** | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub | FAIRDOMHub |

NIB

# Javni podatkovni repozitoriji

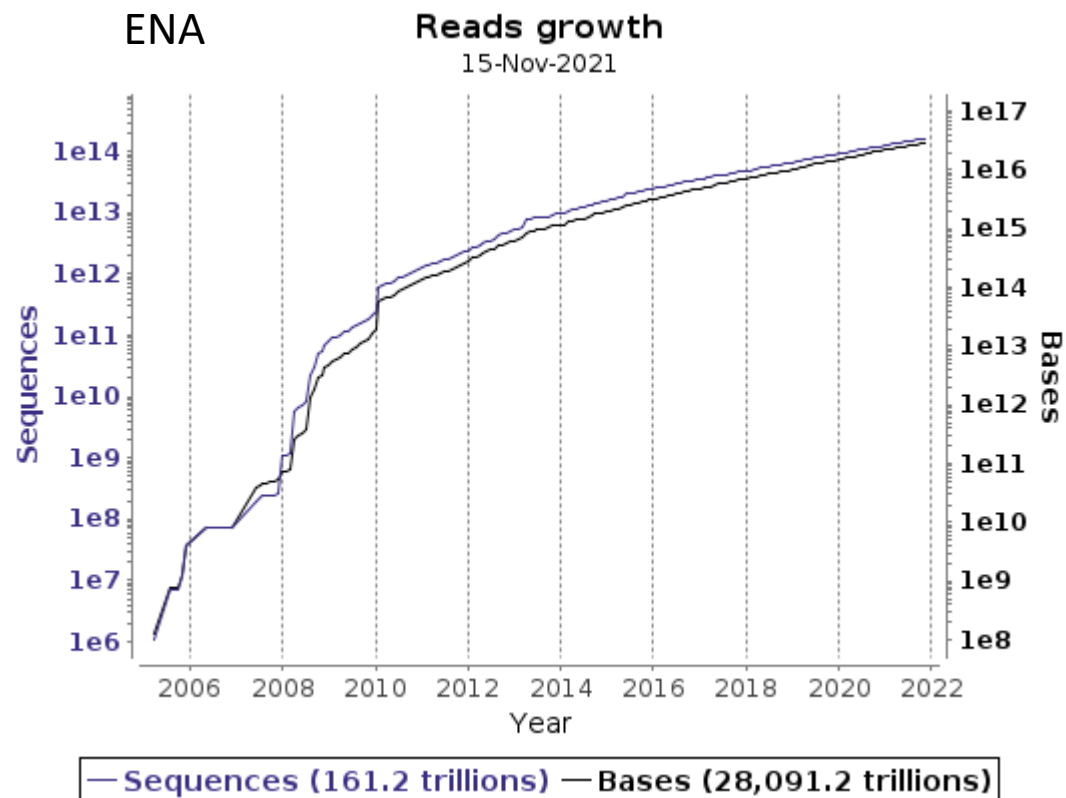**SPECIALIZIRANI REPOZITORIJI ZA SUROVE BIOLOŠKE PODATKE**
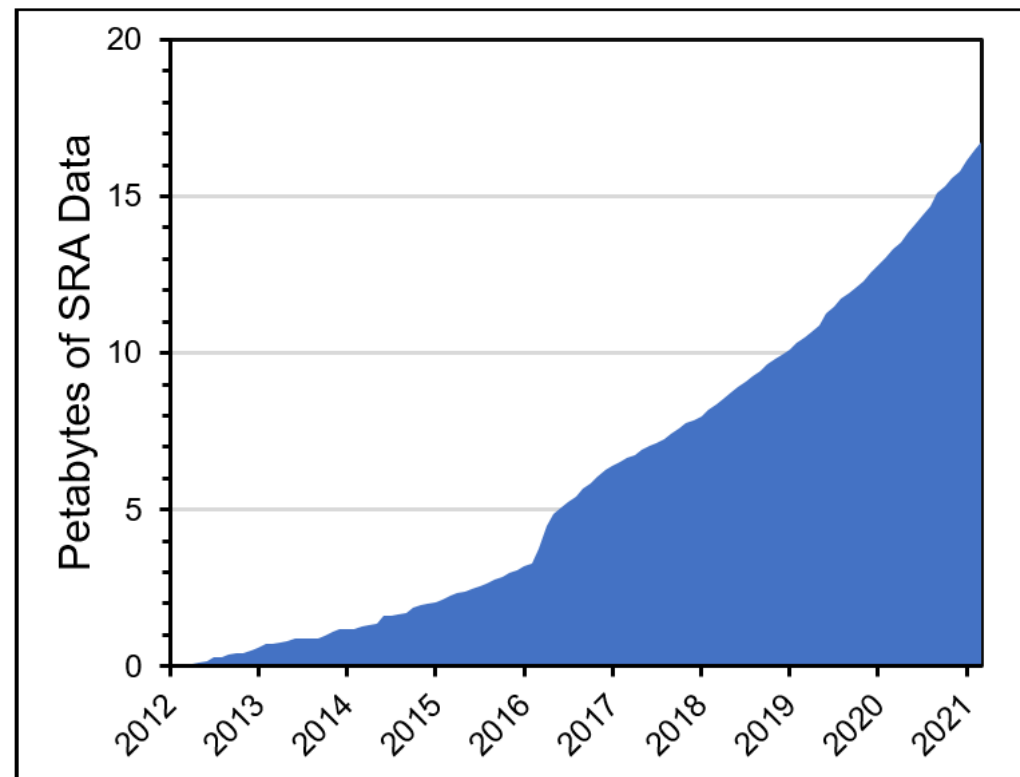
**SPLOŠNI REPOZITORIJI ZA PODATKE**

# Javni podatkovni repozitoriji

**Rast količine podatkov visokozmogljivostnega sekvenciranja DNA (in RNA):**

ENA



slika: EBI ENA
(https://www.ebi.ac.uk/ena/browser/about/statistics)



slika: NCBI Insights
(https://ncbiinsights.ncbi.nlm.nih.gov/2021/08/09/espsss-workshop/#more-6180)

# Deponiranje surovih podatkov v javne repozitorije



https://www.ncbi.nlm.nih.gov/bioproject/PRJNA400633

# Deponiranje surovih podatkov v javne repozitorije

# Kaj pomenijo FAIR podatki za nas?



- enoznačni in perzistentni ID-ji
- strojno berljivi metapodatki

- takojšen dostop
- ali jasna pravila za dostop

- datotečni formati
- ontologije
- genski identifikatorji

- metapodatki z ustreznimi attributi za reanalizo
- licence
- standardi različnih ved

**Kako ravnamo s podatki preden jih deponiramo v javnih repozitorijih?**

# Kako (na NIB) ravnamo s podatki preden jih deponiramo v javnih repozitorijih?

- Surovi podatki
- Metapodatki

- Vmesni in končni rezultati obdelave podatkov in modeliranja

- Publikacije (članki, poročila, predstavitve, posterji etc)

stORK

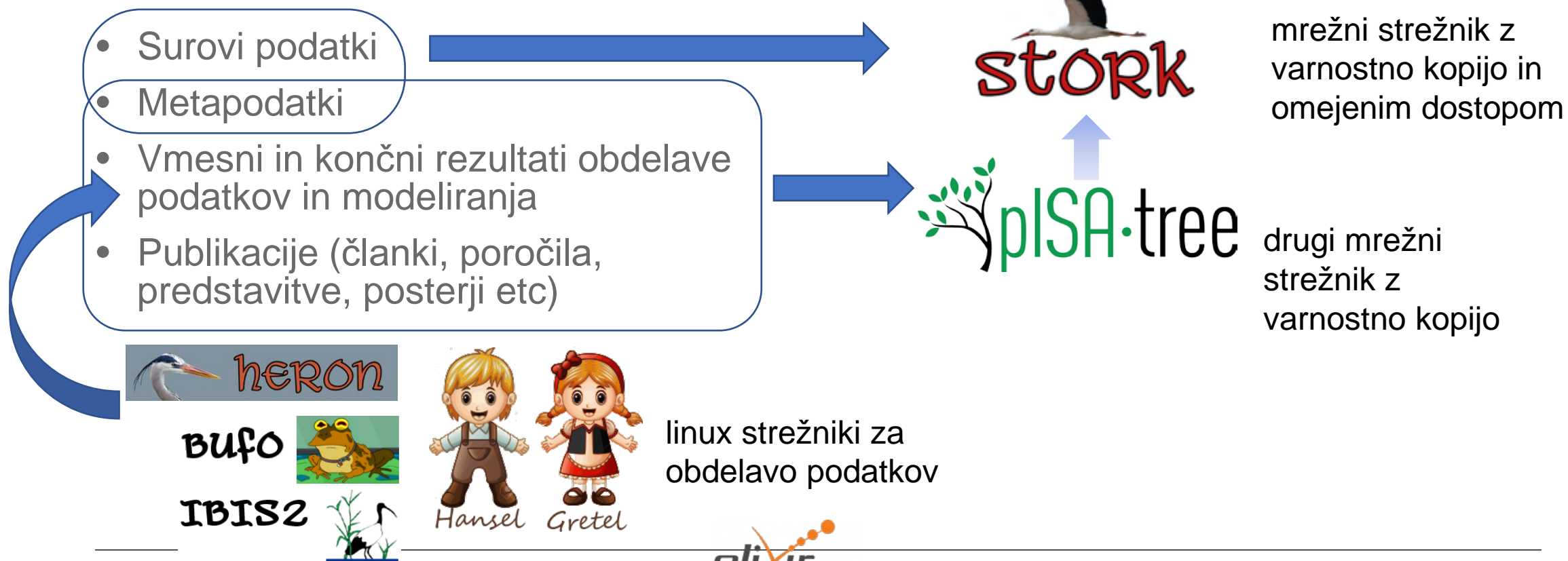mrežni strežnik z varnostno kopijo in omejenim dostopom

# Kako (na NIB) ravnamo s podatki preden jih deponiramo v javnih repozitorijih?

- Surovi podatki
- Metapodatki

- Vmesni in končni rezultati obdelave podatkov in modeliranja

- Publikacije (članki, poročila, predstavitve, posterji etc)

mrežni strežnik z varnostno kopijo in omejenim dostopom

linux strežniki za obdelavo podatkov

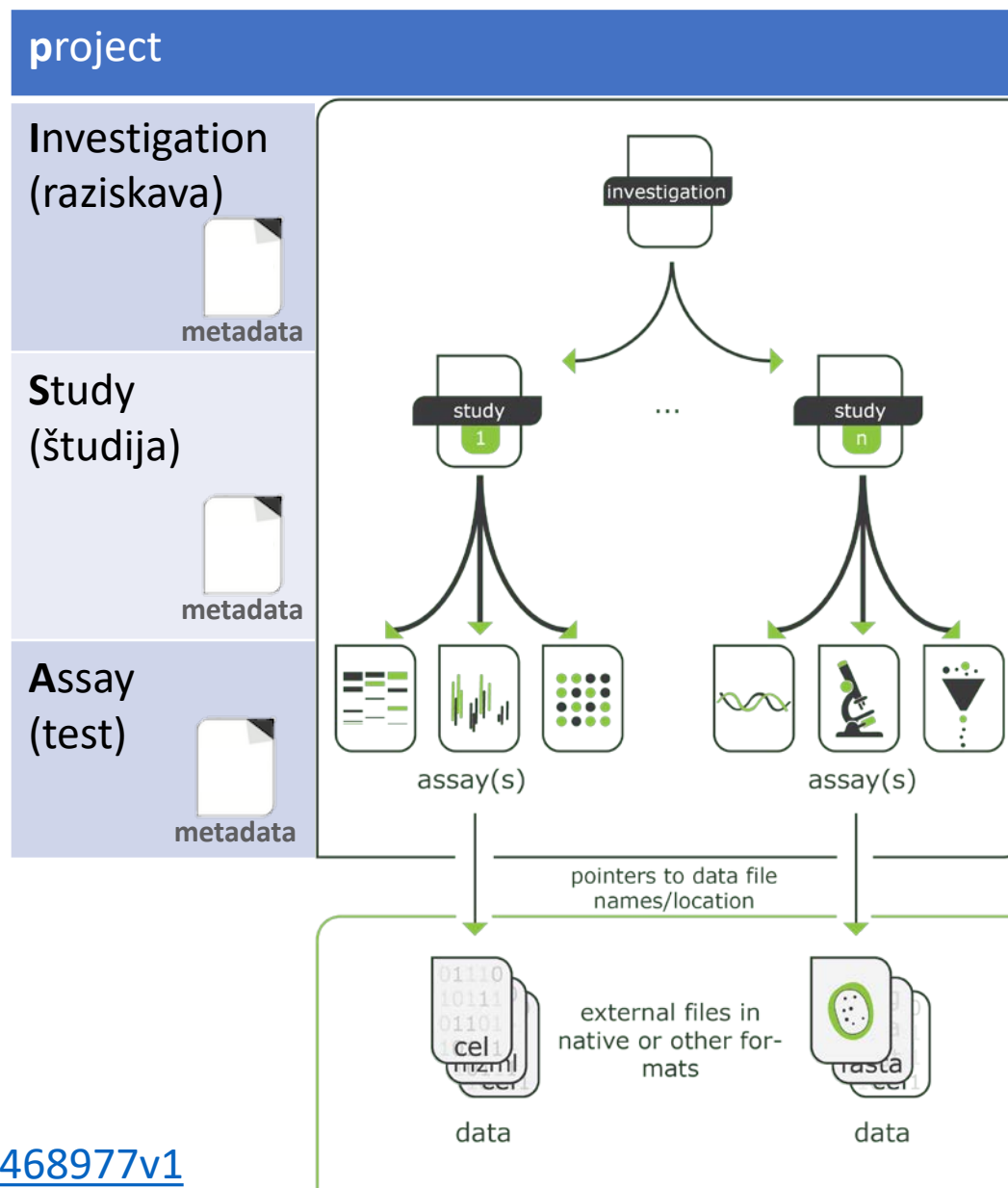# Kako (na NIB) ravnamo s podatki preden jih deponiramo v javnih repozitorijih?

- Surovi podatki
- Metapodatki
- Vmesni in končni rezultati obdelave podatkov in modeliranja
- Publikacije (članki, poročila, predstavitve, posterji etc)

mrežni strežnik z varnostno kopijo in omejenim dostopom

drugi mrežni strežnik z varnostno kopijo

linux strežniki za obdelavo podatkov

# pISA-tree

- Sistem za organizacijo projektnih podatkov (*.**bat** skripte)

- Datotečna drevesna struktura po **ISA specifikacijah**

- metapodatkovne *.txt datoteke v **ISA-Tab format** (sprotno beleženje metapodatkov)

https://github.com/NIB-SI/pISA-tree

Prednatis članka:
https://www.biorxiv.org/content/10.1101/2021.11.18.468977v1



Image: https://isa-tools.org/format/specification.html

# pISA-tree: ustvarjanje nivojev

# pISA-tree: primer lokalne strukture

# pISA-tree: primer lokalne strukture

# pISA-tree: interoperabilnost

# pISA-tree --> FAIRDOMHub.org

- FAIRDOMHub
  - uporablja ISA format
  - na voljo večina licenc za odprto kodo in podatke (Creative Commons, Open Data Commons, …)
  - omogoča pridobitev DOI za podatke

- Prenos v programskem okolju R s paketom *seekr*(https://github.com/NIB-SI/seekr)



https://github.com/NIB-SI/pISA-tree       https://fairdomhub.org/

# pISA-tree --------------------------------------> FAIRDOMHub.org

# ELIXIR Research Data Management Kit

- Spletni vodnik o dobrih praksah ravnanja s podatki za **celoten življenjski cikel podatkov** (know-how, orodja, primeri najboljše prakse)

- Namenjen raziskovalcem, upravljavcem podatkov in oblikovalcem politik

**Your domain**

Plant sciences

Marine metagenomics

Human data

Biomolecular simulation data

Intrinsically disordered proteins

Microbial biotechnology

Epitranscriptome data

Proteomics

Toxicology data

https://rdmkit.elixir-europe.org/

# ELIXIR Research Data Management Kit

# Skupina razvijalcev pISA-tree

**prof. dr. Kristina Gruden**

doc. dr. Špela Baebler

dr. Tjaša Lukan

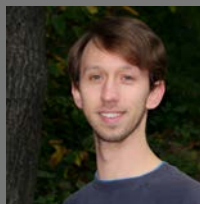prof. dr. Andrej Blejec

dr. Živa Ramšak

Katja Stare

dr. Maja Zagorščak

dr. Anna Coll Rius

Valentina Levak

dr. Marko Petek

marko.petek@nib.si