

The reuse of public datasets in the life sciences: potential risks and rewards

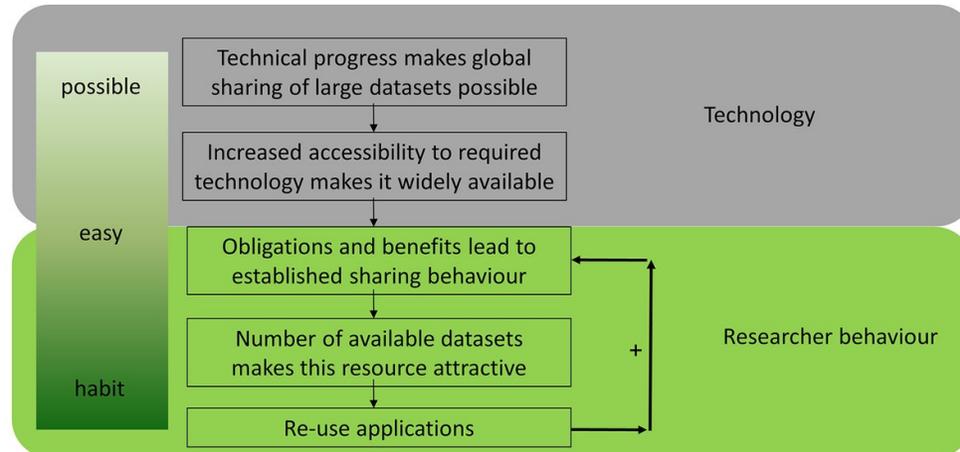
Katharina Sielemann, Alenka Hafner and Boas Pucker

Seminário de Afinidades em Genômica e Bioinformática - Terceiro
encontro de 2024

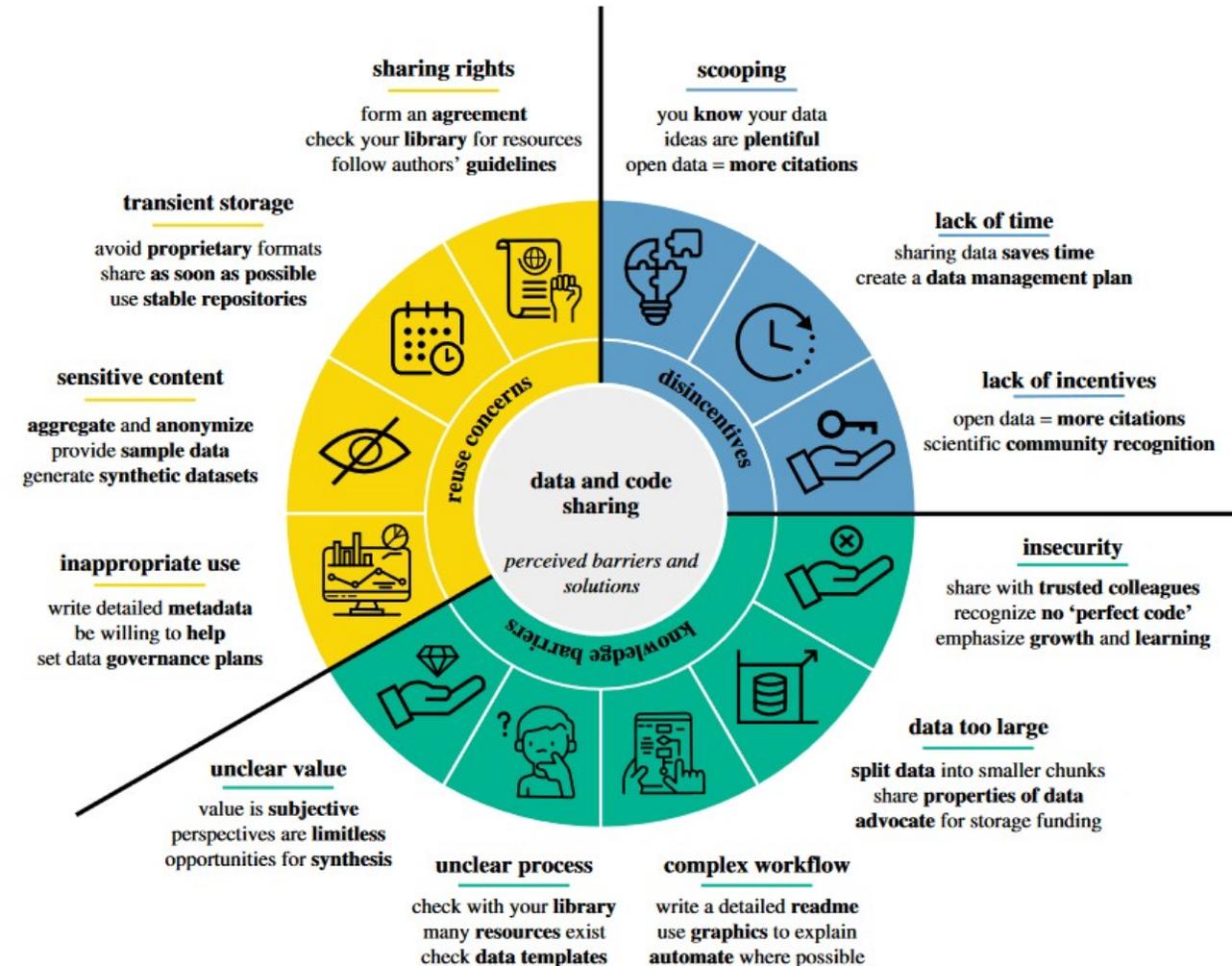
28/08/2024

Big data

- A expansão do “big data” resultou em novas análises para diversos tipos de dados já publicados
- Desenvolvimento de ferramentas de bioinformática para extração de dados já disponíveis:
 - Testes de ferramentas de bioinformática na identificação de padrões
 - Fontes de conjuntos de dados para treinamento de algoritmos de aprendizagem de máquinas



Fonte: Sielemann et al. (2020)

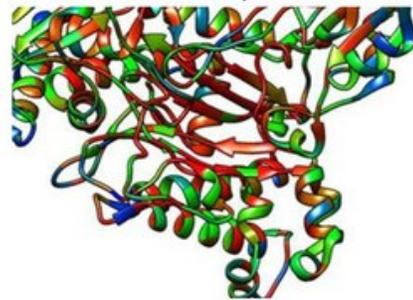


Fonte: Gomes et al. (2022)

Tipos de dados

- Dados primários são diretos e obtidos experimentalmente (sequências, metodologias, grafos, anotações, códigos, medidas, padrões, estruturas, bases de dados primárias)
- Dados derivados são dados de meta-análises e processados (publicações, bases de dados curadas)

```
>seq1
ATCGTTTAGCTAGACCTGATG
ATCCGATCGATTACGTG
>seq2
GACACGATCGTCAGAAATGCA
GTC
>seq3
ACGACAAATCATCTCC
>seq1
MYVRANQEFFK
>seq2
WTSMADCHLKV
>seq3
MGPKLHIGRQEFLIHYWN
NNG
```



```
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') )[-1]
repr_transcripts.update( { repr_trans['id']: transcripts[ repr_trans['id'] ] }
```

 **PDB**e
Protein Data Bank in Europe




PeerJ
Computer Science

JOURNAL FACTSHEET

The multidisciplinary journal for computer science

High-quality, developmental peer review, coupled with industry-leading customer service and an award-winning submission system, means *PeerJ Computer Science* is the optimal choice for your computer science research.

| | |
|--------------------|----------|
| Years publishing | 8 |
| Subjects | 48 |
| Articles published | 1,470 |
| Monthly views | 500,000* |
| Publication speed | 35 days* |
| Impact Factor | 3.8 |
| Scimago Ranking | 0.638 |
| SNIP | 1.094 |
| Citescore | 4.2 |

* All journals
* Median days to 1st decision

Reputation
Publishing high-quality research From the world's top institutions
Professional & experienced team
More than 125 years publishing experience
Industry-leading service
95% of authors recommend us
Institutional publishing
148 institutions signed up

Editors & Advisors
426 Editors and Advisors
Visit Ceff
Google
Valerie Taylor
Texas A&M
Krishnendu Chakrabarty
Duke University
Chieko Asakawa
IBM Japan

Audience
Widely read and cited
500,000 monthly views*
96,000 content alert subscribers
Regular press coverage
In top outlets across US, EU & world
Comprehensively indexed
Web of Science, dbp, CiteSeerX, Google Scholar, Scopus, Pubmed Central, DOAJ, ++

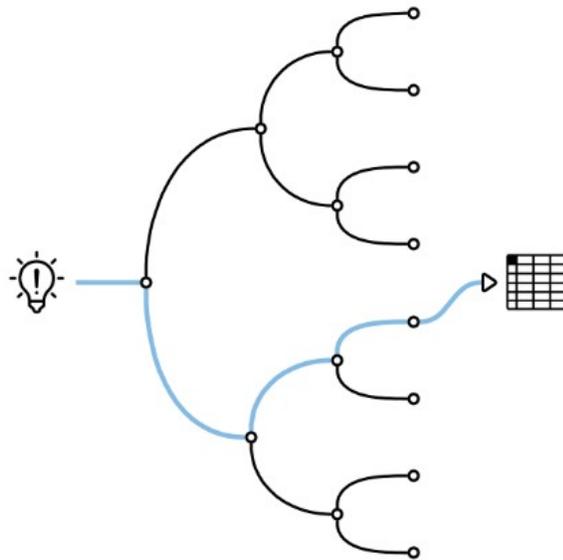
Publishing Made Easy
Relaxed reference formatting
Clarity is the only requirement
Press release & sharing tools
Maximize readership of your work
Great looking articles
Fully typeset article proofs
Straightforward submission
Easy to use, fast and loved by authors

Peer review
Robust, Developmental review
1 editor, 2+ reviewers
Excellent review quality
Above average quality and depth
Signed reviews (optional)
Opt in to fair and transparent science

peerj.com/computer-science

Potenciais

- O estabelecimento de normas de compartilhamento de dados devem ser introduzidas por obrigações nas publicações, incluindo descrição e padronização



Fonte: Borghi e Van Gulick. (2022)

| | Project Planning | Data Collection and Analysis | Data Publication and Sharing |
|------------------------------------|--|--|---|
| Data Management Activities | <ul style="list-style-type: none"> • Data Management Planning | <ul style="list-style-type: none"> • Saving and backing up files • Organizing files • Formatting and describing data according to standards • Maintaining documentation and metadata | <ul style="list-style-type: none"> • Preserving data and other materials (e.g. reagents, code) • Assigning persistent identifiers |
| Open Science Activities | <ul style="list-style-type: none"> • Planning for open (e.g. including data sharing in consent forms) | <ul style="list-style-type: none"> • Using open source tools • Using transparent methods and protocols | <ul style="list-style-type: none"> • Sharing data • Publishing research reports openly (e.g. open access publishing). |
| Reproducibility-Related Activities | <ul style="list-style-type: none"> • Preregistering study aims and methods • Using appropriate research designs (including sufficient statistical power) | <ul style="list-style-type: none"> • Preventing methodological issues (e.g. p-hacking, HARKing) • Implementing quality control measures | <ul style="list-style-type: none"> • Preventing publication bias • Following reporting guidelines |

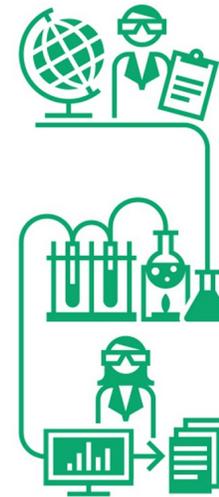
Fonte: Borghi e Van Gulick. (2022)

Potenciais

- A padronização exige precauções:
 - A integridade da fonte deve ser verificada
 - Possíveis vieses devem ser analisados
 - Todas as informações relevantes devem estar disponíveis (metadados)
 - A integração de dados só pode ser realizada com dados semelhantes
 - A qualidade deve ser alta (Phred scores, comprimento e pareamento de sequências, contig/scaffold N50, etc)

| Name | Explanation | Prevention |
|--------------------|---|---|
| Allocation bias | Systematic difference in the assignment of participants to the treatment and control group in a clinical trial. For example, the investigator knows or can predict which intervention the next eligible patient is supposed to receive due to poorly concealed randomization. | - Randomization with allocation concealment |
| Attrition bias | Attrition occurs when participants leave during a study that aims to explore the effect of continuous exposure (dropouts or withdrawal). For example, more dropouts of patients randomized to an aggressive cancer treatment. | - Good investigator-patient communication - Accessibility of clinics - Incentives to continue |
| Confounding bias | An artificial association between an exposure and an outcome because another variable is related to both the exposure and outcome. For example, lung cancer risk in coffee drinkers is evaluated, ignoring smoking status (smoking is associated with both coffee drinking and cancer). A challenge is that many confounders are unknown and/or not measured. | - Randomization (can address unmeasured confounders) When randomization is not possible: - Restriction to one level of the confounder - Matching on the levels of the confounder - Stratification and analysis within strata - Propensity score matching |
| Immortal time bias | Survival beyond a certain time point is necessary in order to be exposed (participants are "immortal" in that time period). For example, discharged patients are analyzed but were included in the treatment group only if they filled a prescription for a drug 90 days after discharge from hospital. | - Group assignment at time zero - Time-dependent analysis may be used |
| Information bias | Bias that arises from systematic differences in the collection, recall, recording, or handling of information. For example, blood pressure in the treatment arm is measured in the morning and for the control arm in the evening. | - Standardized data collection - Data collection independent from exposure or outcome (e.g., by blinding of intervention status/exposure) - Use of objective measurements |
| Publication bias | Occurs when only studies with a positive or negative result are published. Affects meta-analyses from systematic reviews and harms evidence-based medicine | - Writing a study protocol and preregistration - Publishing study protocol or registered report - Following reporting guidelines |

Fonte: Schwab et al. (2022)



Planning

1. Specify your research question
2. Write and register a study protocol
3. Justify your sample size
4. Write a data management plan
5. Reduce bias

Execution

6. Avoid questionable research practices
7. Be cautious with interpretations of statistical significance
8. Make your research open

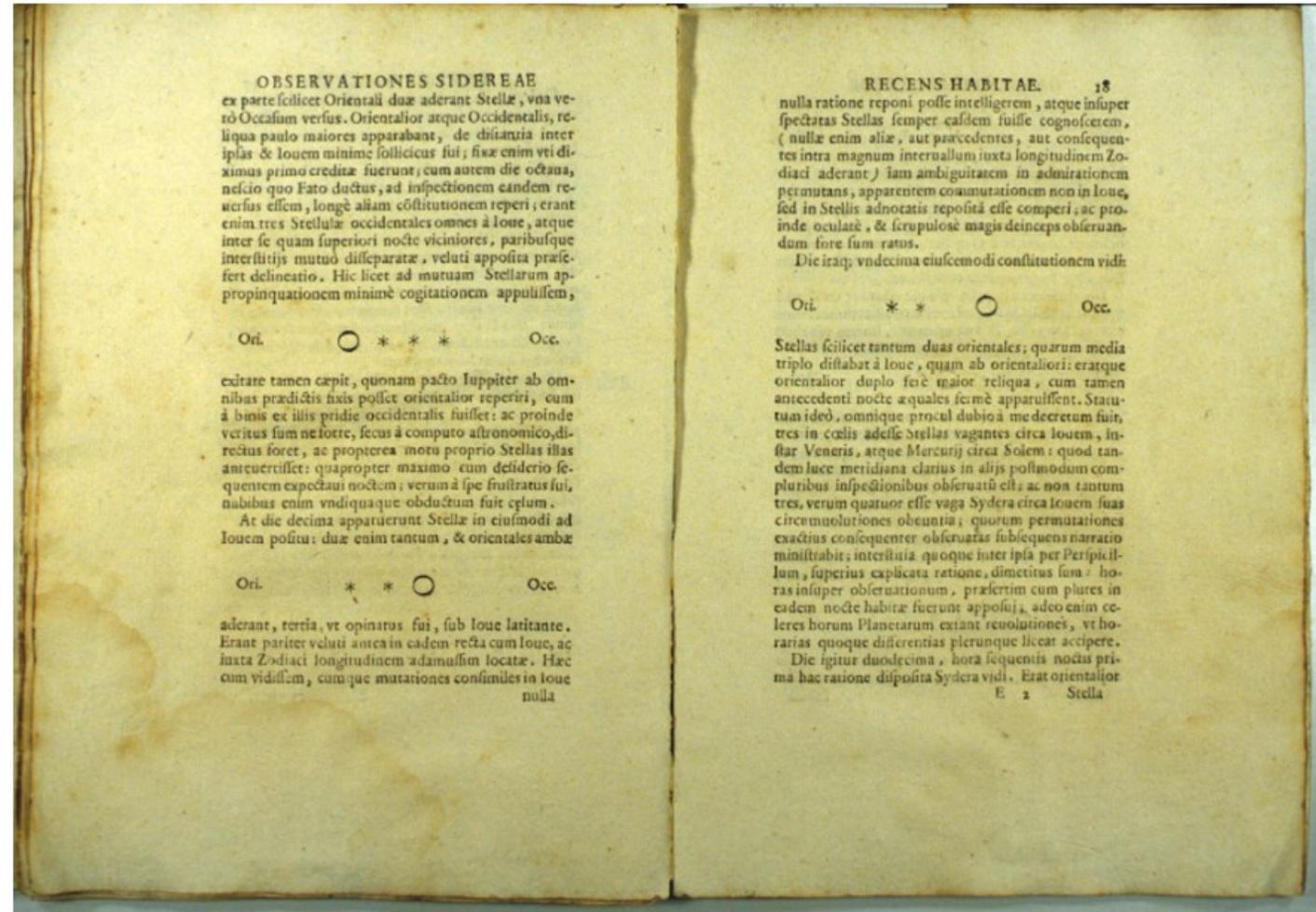
Reporting

9. Report all findings
10. Follow reporting guidelines

Fonte: Schwab et al. (2022)

Potenciais

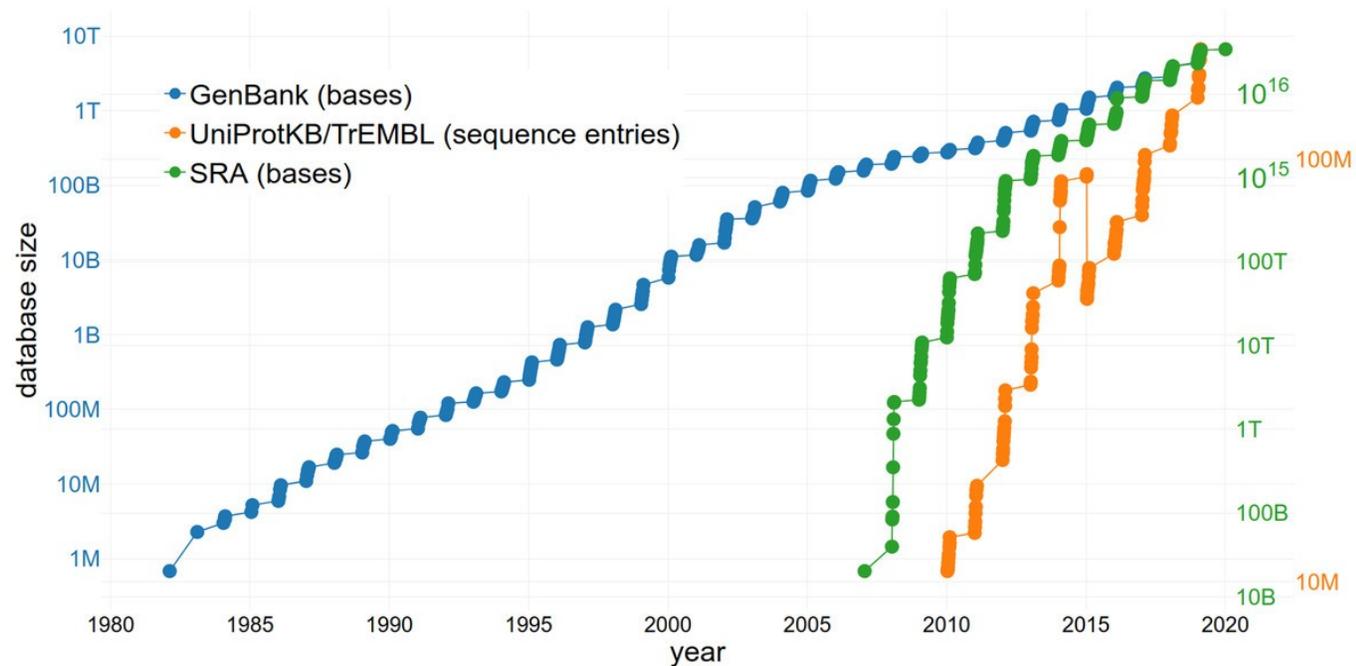
- A manutenção de dados favorece também a integração de dados de diferentes fontes:
 - Fornecem análises necessárias para hipóteses que não podem ser testadas por conjuntos únicos de dados
 - A reutilização de dados supera limitações financeiras de muitos grupos menores para o delineamento de experimentos complexos, uma vez que a geração de novos dados pode ser custosa e laborosa



Fonte: Goodman et al.
(2014)

Potenciais

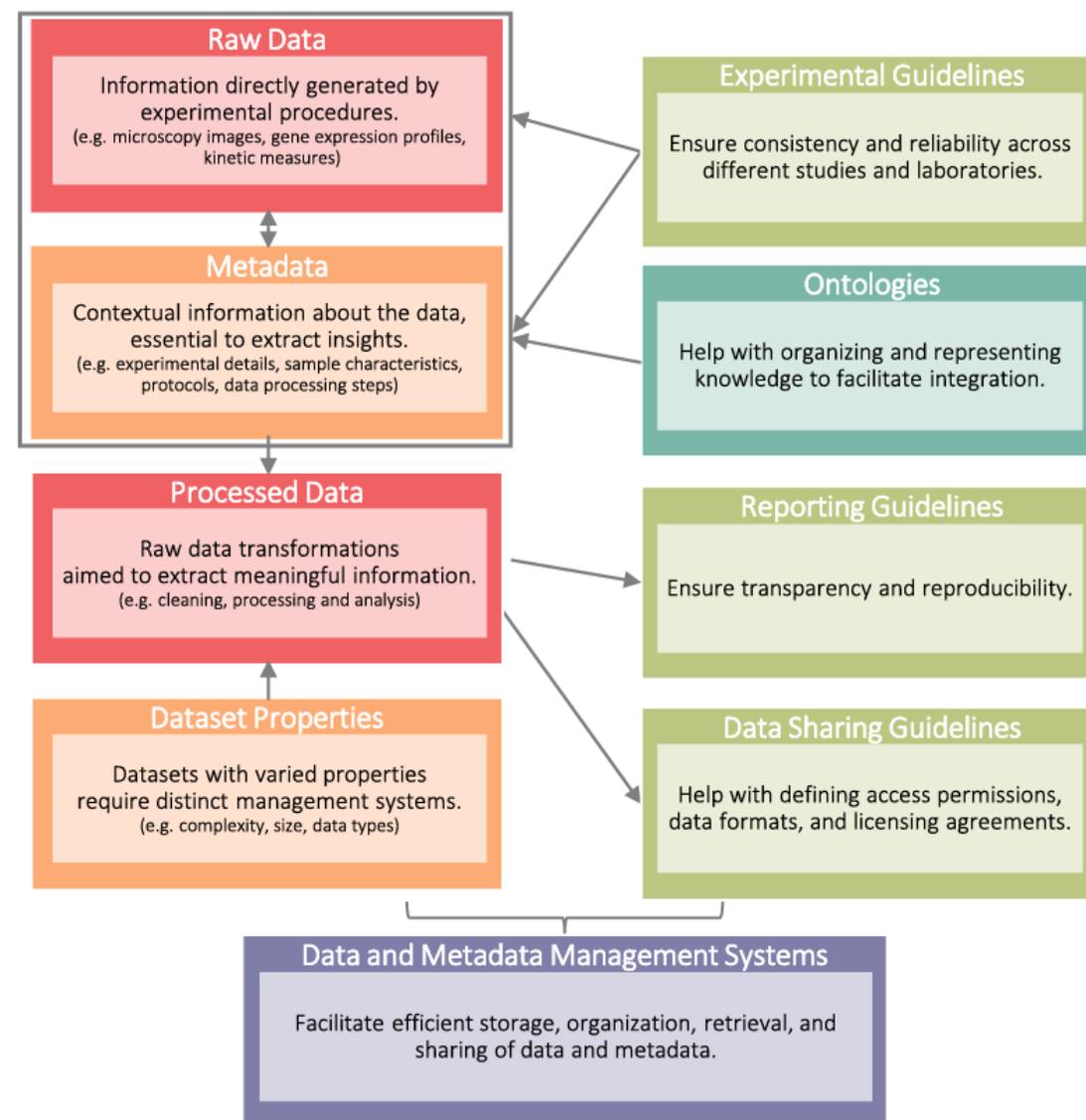
- A disponibilização pública dos dados previne a perda de informação e a geração de conjuntos redundantes e/ou duplicados:
 - Manutenção concisa das bases de dados
 - 20% das publicações de metagenomas entre 2016 e 2019 não são acessíveis ao público



Fonte: Sielemann et al.
(2020)

Potenciais

- A reprodutibilidade da reutilização de dados pode ser melhorada pela padronização, documentação e organização de fluxos de trabalho:
 - A integração de diferentes fontes envolve a comparação entre fluxos de trabalho distintos, variáveis indevidamente registradas, apresentações não unificadas, e falta de disponibilidade de dados brutos
 - Exige uma curagem manual
 - Padronização do formato de linguagem e arquivos (XML, FASTA, FASTQ, SAM/BAM)



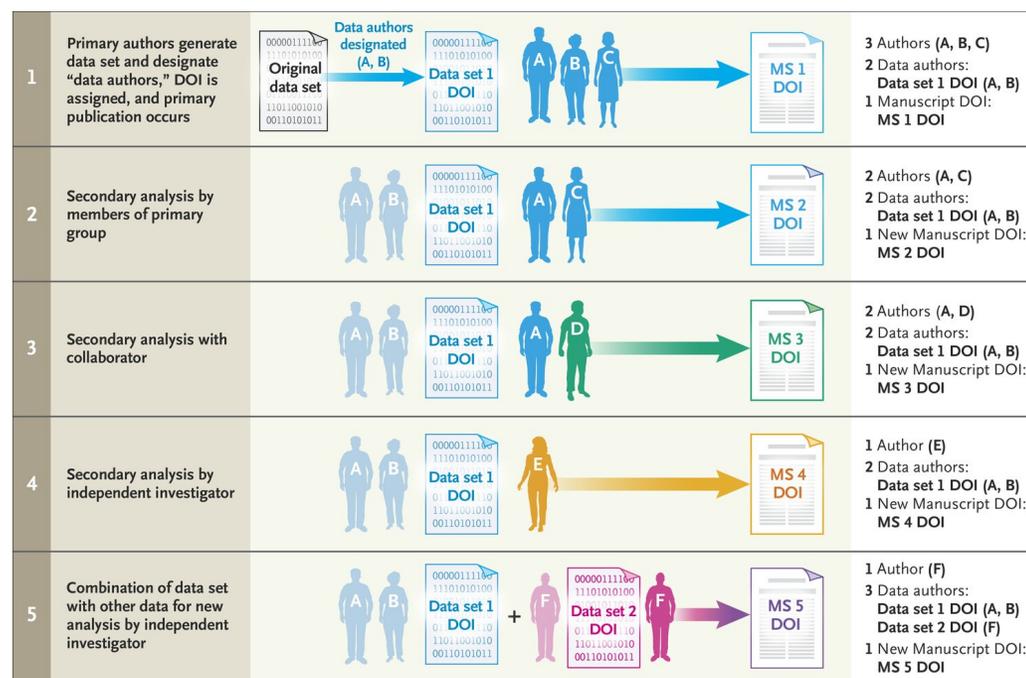
Potenciais

- Existem guias para a padronização do compartilhamento de metadados:
 - Findable Accessible Interoperable Reusable (FAIR)
 - The Transparency and Openness Promotion (TOP)
 - Open Data in a Big Data World (Open Data in a Big Data World, 2016)
 - Beijing Declaration

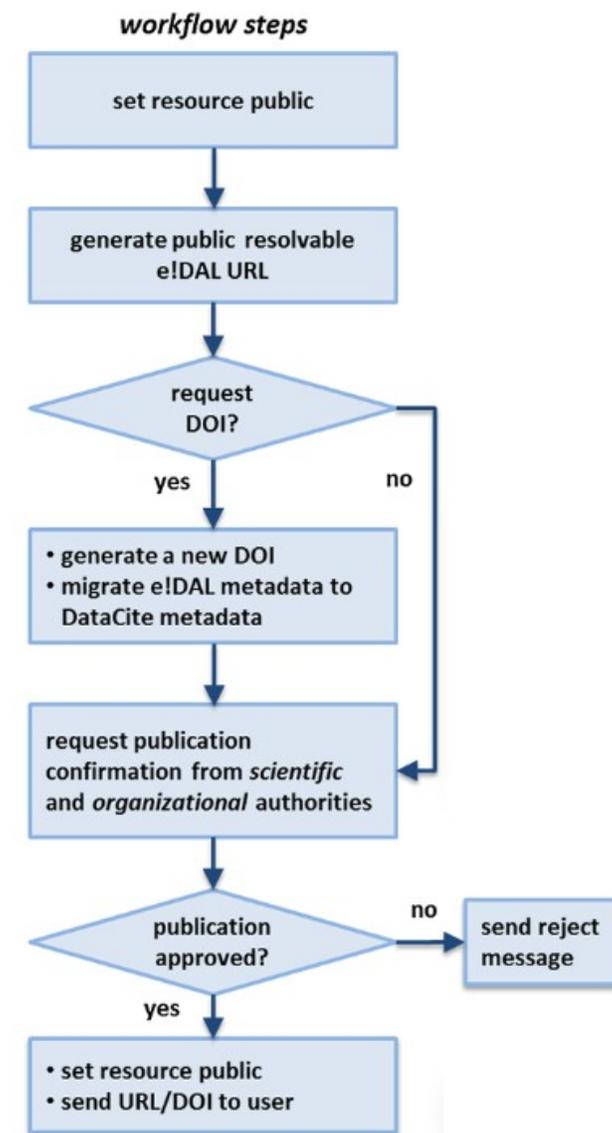
| FAIR Principle | For Infrastructure | For Researchers |
|----------------------|---|---|
| Findable | <p>Data and metadata should be easy to find by both humans and computers.</p> <p>In practice, this means that data should be assigned unique and persistent identifiers, described using rich metadata, and registered or indexed in a searchable resource.</p> | <p>Research teams should implement standardized practices related to organizing files (e.g. standardized file naming conventions) so data can be found when needed.</p> <p>When data is made available to others, it should- whenever possible- be uploaded to a repository that assigns a persistent identifier (e.g. DOIs, RRIDs, etc) and describes datasets with standardized metadata.</p> <p>Complete and high quality metadata should be added so data can be discovered and linked to related resources (e.g. related paper DOIs, author ORCIDs).</p> |
| Accessible | <p>There is a clearly defined method for accessing the data. Data should be retrievable by its identifier using a standardized protocol that is open, free, and universally implementable. Metadata should be accessible even when data is no longer available.</p> | <p>Data is available through a clearly defined process. Members of the research team should be able to access raw data, intermediate products, and other research materials.</p> <p>When data and other materials are made available to others, there should be a clear path to gaining access. The terms by which the data will be made available (e.g. to whom, when, and for what purpose) should be articulated and abided by.</p> |
| Interoperable | <p>Data should be usable across a range of applications and workflows. Data should use formal, accessible, shared, and broadly applicable models for knowledge representation.</p> | <p>Data should be structured in a standard way so it can be easily combined with other similarly structured datasets. In practice, this means implementing a range of practices such as describing and organizing data (e.g. applying appropriate metadata, maintaining data dictionaries) and saving files in open or non-proprietary file formats.</p> |
| Reusable | <p>Metadata and data should be described following relevant community standards and have clearly defined conditions for reuse including a machine-readable license</p> | <p>Data should be saved, organized, and described with its future (re)use in mind. A future user may be a member of the research team who is returning to the data after a period of several months or years or another researcher who is (re)using the data for another purpose.</p> |

Potenciais

- Autores podem ser beneficiados pelo fornecimento de dados de alta qualidade e com boa anotação:
 - A prática deve fomentar as citações e reconhecimentos adequados
 - Podem resultar em citações adicionais para um mesmo conjunto de dados



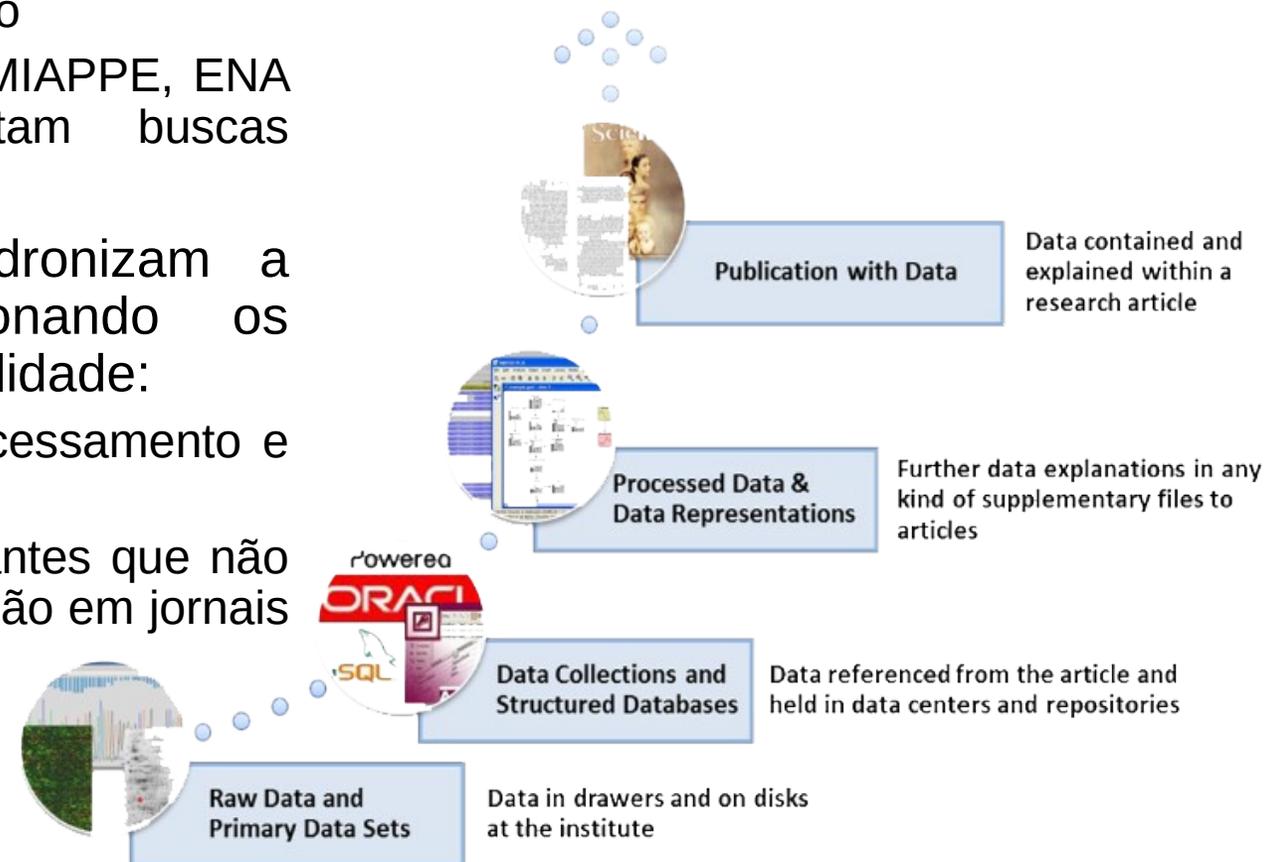
Fonte: Bierer et al. (2017)



Fonte: Arend et al. (2014)

Potenciais

- A padronização exige uma curagem manual:
 - Cada projeto apresenta vocabulário próprio
 - Padronização de metadados (MassIVE, MIAPPE, ENA e SRA) e de vocabulário facilitam buscas automatizadas
- Jornais de dados (*data papers*) padronizam a submissão dos metadados, solucionando os problemas associados ao controle de qualidade:
 - Descrição dos métodos para coleta, processamento e verificação dos dados
 - Permitem a publicação de dados importantes que não são suficientes para compor uma publicação em jornais comuns.



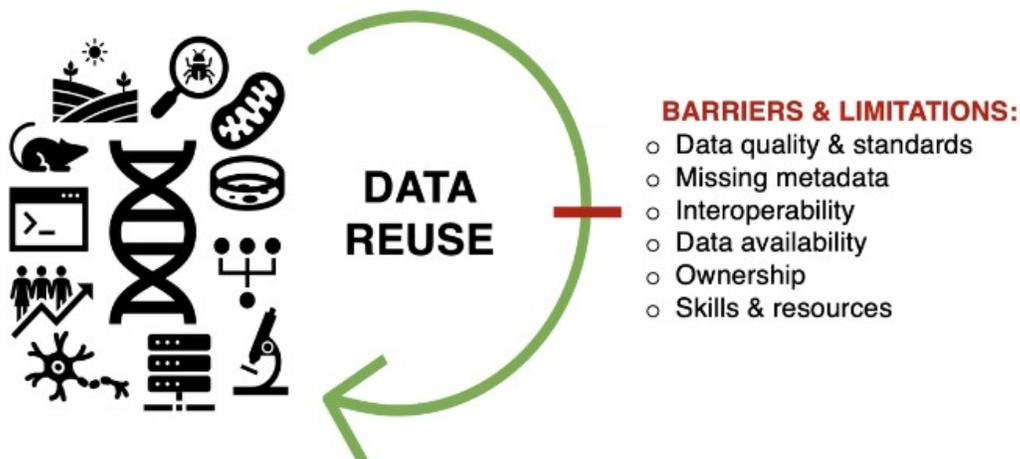
Ética

- Direitos de uso:
 - Reuso justo (para novos propósitos)
 - Reprodução de estudos prévios com dados publicados
 - Reuso injusto (plágio)
- Parasitismo de pesquisa:
 - A reutilização deve gerar novos resultados, e não apenas impulsionar publicações anteriores.
 - *The Research Parasite Awards* premia autores que adequadamente reutilizaram dados



Limitações

- A revisão por pares geralmente não envolve a avaliação dos conjuntos de dados e sua descrição:
 - A qualidade dos dados é questionável quando são de responsabilidade pública
 - Metodologia e design experimental incompletas
 - Propagação de ruídos nos dados com a reutilização
 - Os conjuntos de dados não são normalizados, impedindo a integração



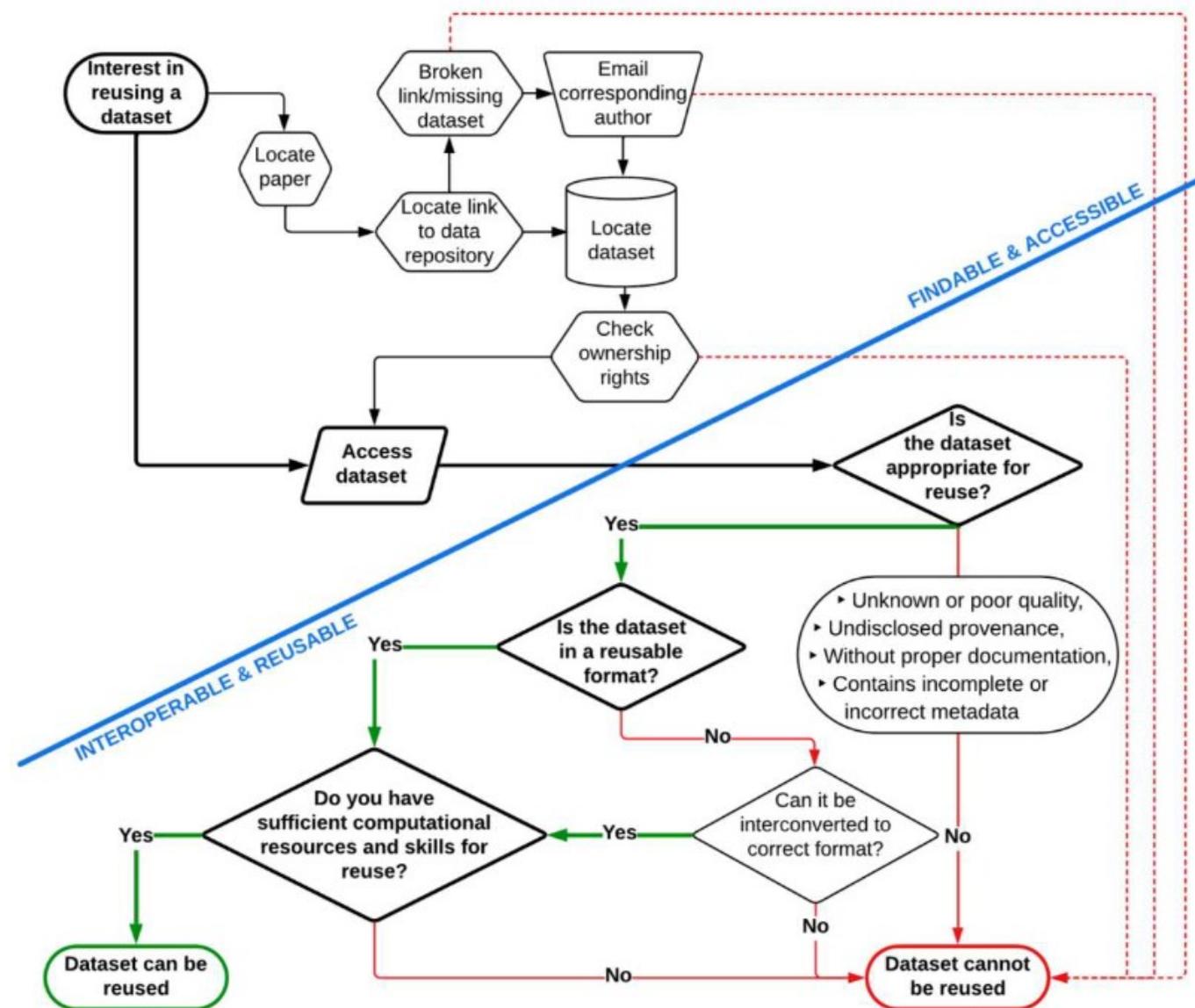
Fonte: Hafner et al. (2024)



Fonte: Raman et al. (2021)

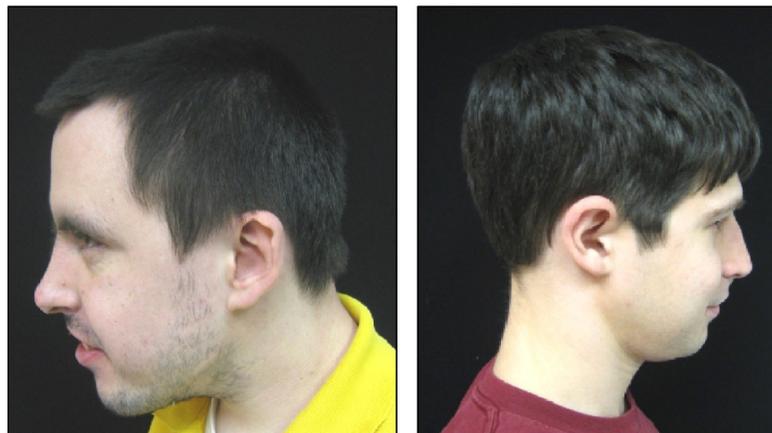
Limitações

- Sequências podem ser depositadas com identificadores diferentes, impedindo sua integração e escondendo possíveis informações.
- A anotação não é padronizada, possibilitando a inserção de um mesmo conjunto de dados com diferentes identificações
- A montagem de genomas baseada em dados públicos podem ser prejudicadas pela presença de contaminações não identificadas
- Falsos positivos em experimentos podem não ser identificados em dados públicos

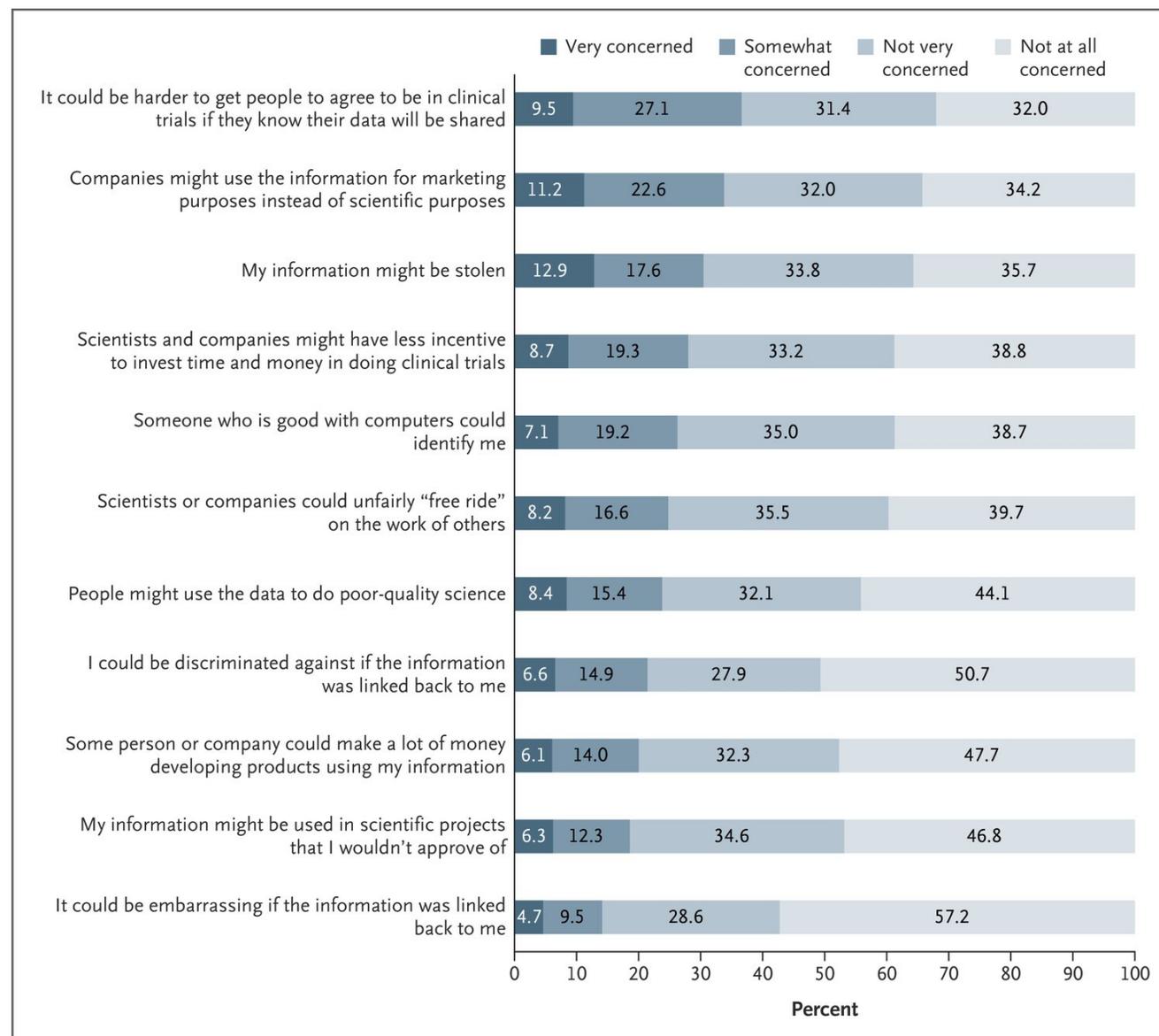


Limitações

- Dados clínicos impactam os pacientes envolvidos e devem ser devidamente assegurados para reutilização
 - Pacientes não devem ser identificados
 - Dados genômicos (risco de privacidade para a progênie)



Fonte: McDonald-McGinn et al. (2011)



Fonte: Mello et al. (2018)

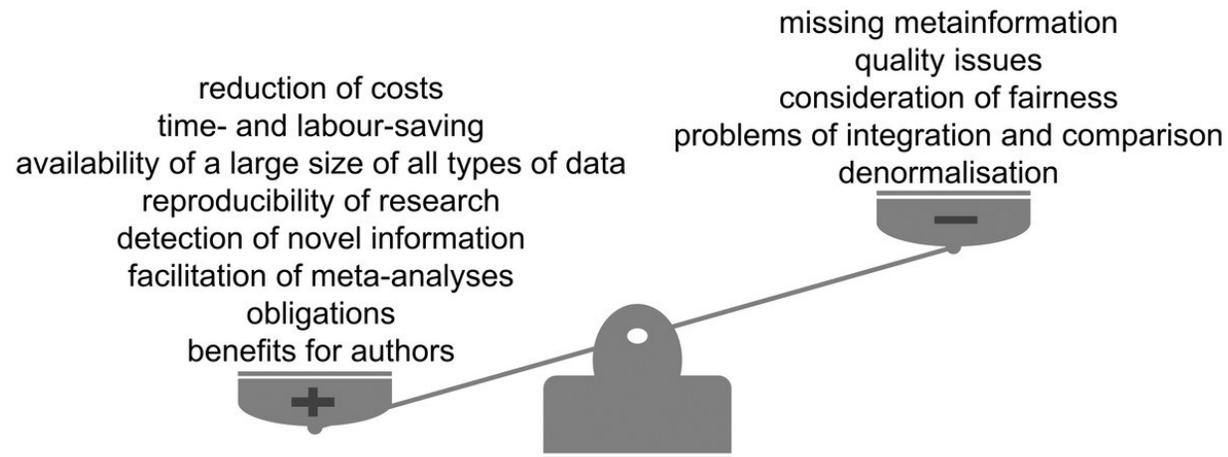
Balanço

■ Benefícios:

- Redução de custos
- Redução de trabalho manual
- Disponibilidade de diversos tipos de dados
- Reprodutibilidade de trabalhos
- Detecção de novas informações
- Facilidade em meta análises
- Benefícios aos autores.

■ Precauções:

- Perda de informação
- Perda de qualidade
- Uso justo de dados
- Problemas de integração e comparação
- Falta de normalização

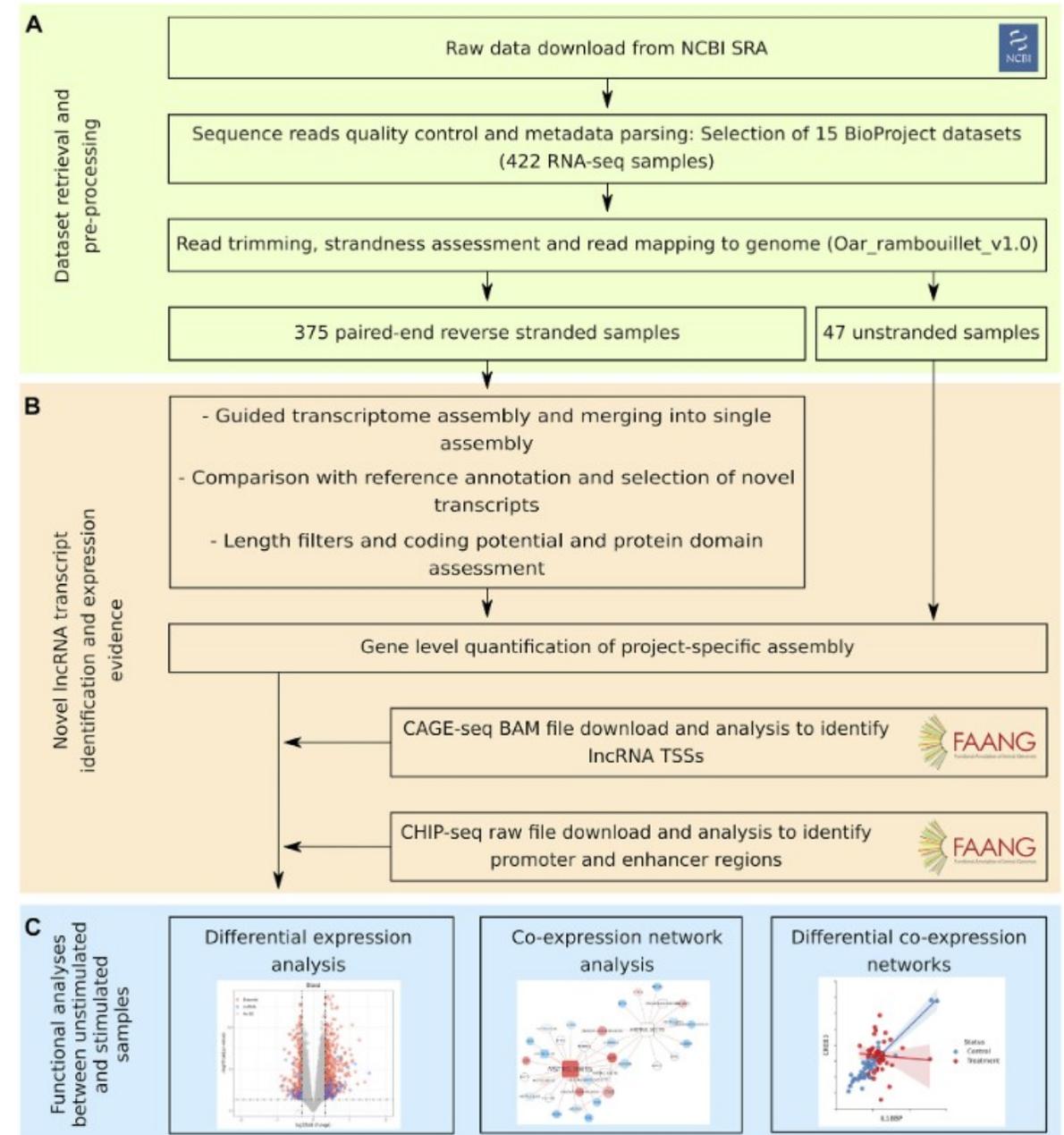


Examples

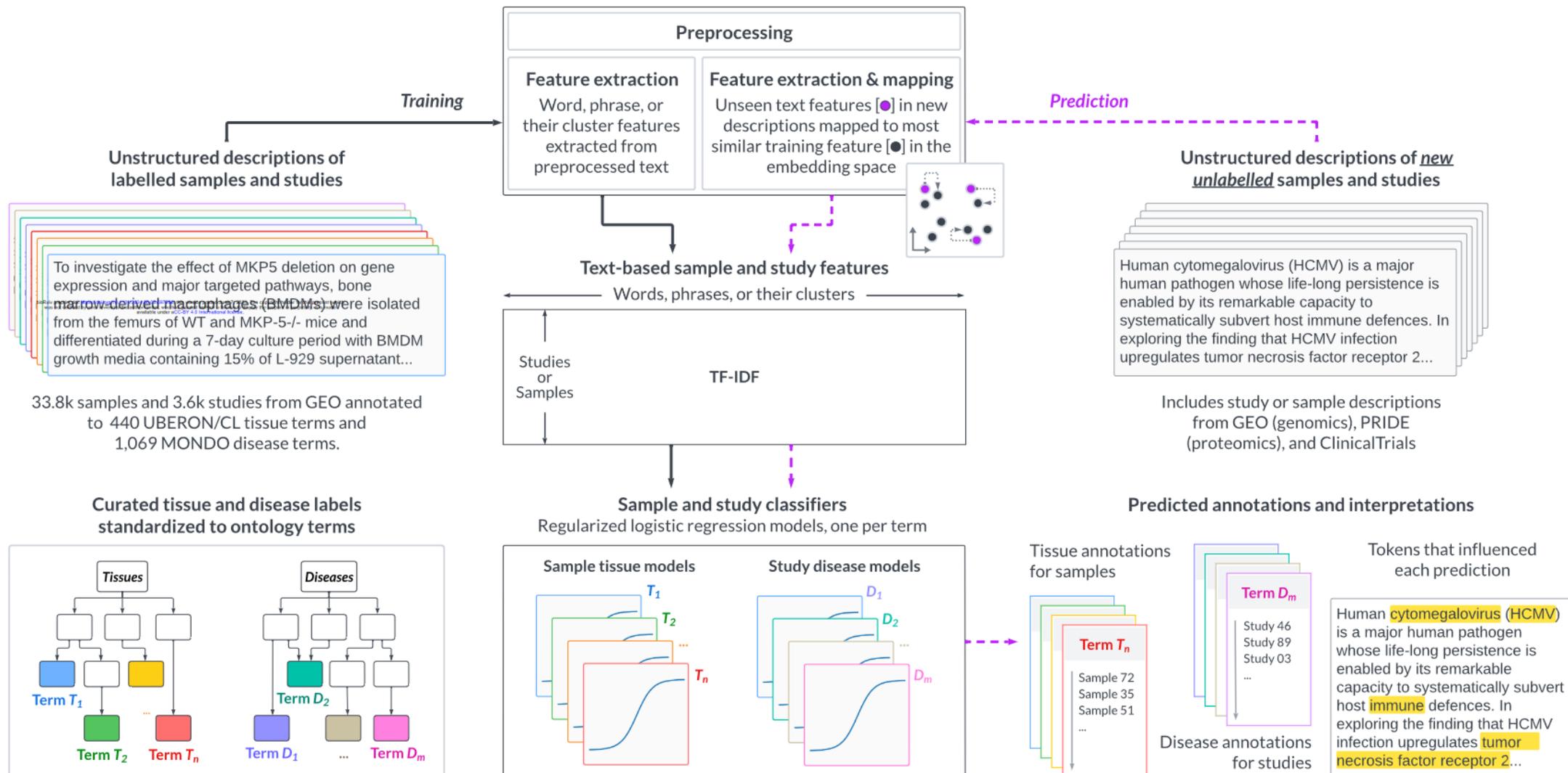
| FMD | Title | First Author | Journal/Book | Publication Year | DOI | 2nd Author | Journal/Book | Publication Year | DOI | 3rd Author | Journal/Book | Publication Year | DOI |
|----------|--|----------------|------------------------------|------------------|---------------------------------|------------|---|------------------|------------------------------|---------------------|---------------------------------------|------------------|------------------------------|
| 2848073 | Clinical Data Base or Secondary Use: Current Status and Potential Future Progress | Meyreite SM | Year's Med Inform | 2017 | 10.1093/imj/njz007-007 | 2947920 | A Web-based Tool to Enhance Monitoring and Attention in Care for Tuberculosis Affected Patients | 2017 | 10.1186/s12859-016-1131-1 | suaimin B | Stout Health Technol Inform | 2017 | |
| 2604621 | Intelligence-based data warehouse environment to enable the reuse of electronic health record data | Martini-Hell J | Int J Med Inform | 2017 | 10.1016/j.ijmedinf.2016.05.016 | 2784608 | GENVA: aggregation and analysis of gene expression signatures from related studies | 2016 | 10.1186/s12859-016-1131-1 | Gundersen GW | BMC Bioinformatics | 2016 | 10.1186/s12859-016-1131-1 |
| 3026916 | Data Reuse Through Anesthesia Data Warehouse: Searching for New Use Contexts | Lamer A | Stud Health Technol Inform | 2018 | | 23952476 | Electronic health records: new opportunities for clinical research | 2017 | 10.1111/imj.12119 | Corevrets P | J Intern Med | 2017 | 10.1111/imj.12119 |
| 2828186 | Methods and norms affecting scientific data reuse | Canoy R | PLoS One | 2017 | 10.1371/journal.pone.0182426 | 2701470 | Data Extraction and Management in Networks of Observational Health Care Databases for Scientific Research: A Comparison of EU-ADR, OMOP, Mini-Sentinel and MATRIC Strategies | 2017 | 10.13683/2327-9214.1189 | Gini R | ECEMS (Wash DC) | 2016 | 10.13683/2327-9214.1189 |
| 3030551 | Changes in Data Sharing and Data Access Practices and Perceptions among Scientists Worldwide | Thompson C | PLoS One | 2018 | 10.1371/journal.pone.0219420 | 30251936 | An open-HDL Approach to Detailed Clinical Model Development: Tobacco Smoking Summary Architecture as a Case Study | 2019 | 10.1093/ibd/ibz009-1603104 | Wier PC | Appl Clin Inform | 2019 | 10.1093/ibd/ibz009-1603104 |
| 3000305 | Update on Data Reuse in Health Care | Sahlic C | Year's Med Inform | 2017 | 10.1093/imj/njz007-013 | 29125569 | From record keeping to scientific research: obstacles and opportunities for research with electronic health records | 2017 | 29125569 | Schotte RA | Ned Tijdschr Geneesk | 2017 | |
| 3134773 | Credit card generation for data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 29151881 | Perspectives for medical informatics: revisiting the electronic medical record for clinical research | 2017 | | Prochaski HJ | Methods Inf Med | 2009 | |
| 3130959 | Data reuse and the open data citation advantage | Flaxman HA | PLoS One | 2013 | 10.1371/journal.pone.0101175 | 31309550 | Temporal variability analysis reveals biases in electronic health records used at hospital periods: reengineering interventions over seven years | 2019 | 10.1371/journal.pone.0220969 | PiVincuz-Bentini FJ | PLoS One | 2019 | 10.1371/journal.pone.0220969 |
| 2926120 | Review of adults worldwide who reuse of health data among people in the European Union: The primary of purpose and the common good | Stangorill L | Health Policy | 2019 | 10.1016/j.healthpol.2019.03.012 | 27322245 | Testing of Triggers by Data Mining of Epilepsy Patients' Structured Nursing Records | 2018 | 10.1006/j.jmig.2018.01.018 | Kinnunen UM | Stud Health Technol Inform | 2016 | |
| 2819643 | Reuse of public genome-wide gene expression data | Rung J | Nat Rev Genet | 2013 | 10.1038/nrg3104 | 23304237 | ClinData Express—a metadata driven clinical data management system for secondary use of clinical data | 2012 | 10.1038/nrn.2016.08.004 | Li Z | AMIA Annu Symp Proc | 2012 | |
| 1412480 | Clinical trial data reuse—overcoming complexities in trial design and data sharing | Wilkinson T | Trials | 2010 | 10.1186/1745-2875-11-2672-6 | 28139283 | Reusability of coded data in the primary care electronic medical record: a dynamic context concerning cancer diagnoses | 2017 | 10.1038/nrn.2016.08.004 | Sollie A | Int J Med Inform | 2017 | 10.1038/nrn.2016.08.004 |
| 3002178 | Intelligence-based data warehouse environment to enable the reuse of electronic health record data | Canoy R | PLoS One | 2017 | 10.1371/journal.pone.0182426 | 23633602 | Integrated database of information from structural genomics experiments | 2015 | 10.1107/2090444913007128 | Aadaa Y | Acta Crystallogr D Biol Crystallogr | 2013 | 10.1107/2090444913007128 |
| 3073930 | Proteomics data reuse with MAUD-IO | Duenen A | Bioinformatics | 2019 | 10.1093/bioinformatics/bty408 | 27529152 | Quantitative monitoring of Arabidopsis thaliana growth and development using high-throughput plant phenotyping | 2016 | 10.1038/nature.2016.55 | Arénd D | Sci Data | 2016 | 10.1038/nature.2016.55 |
| 2550889 | Carotid-Flow for CVD and Outcomes data reuse | Jan H | Bioinformatics | 2018 | 10.1093/bioinformatics/bty023 | 24747879 | The need for harmonized structured documentation and chains of secondary use - results of a systematic analysis with automated form comparison for prostate and breast cancer | 2014 | 10.1016/j.jmig.2014.04.008 | Krumm R | J Biomed Inform | 2014 | 10.1016/j.jmig.2014.04.008 |
| 2918814 | A Data Quality Assessment Guideline for Health Record Data Reuse: A Conceptual Best Practice Framework and Procedure Model | Woolcock NG | ECEMS (Wash DC) | 2017 | 10.5334/egen.218 | 23244553 | Semantically enabling a genome-wide association study database | 2014 | 2014.1186/2041-1480-3-9 | Beck T | J Biomed Semantics | 2014 | 2014.1186/2041-1480-3-9 |
| 2676124 | SHOT: Systematic Planning of Intelligent Search of Integrated Clinical Routine Data: A Conceptual Best Practice Framework and Procedure Model | Huall WO | Methods Inf Med | 2010 | 10.1016/j.mim.2010.03.005 | 24048520 | Clinical professional governance for detailed clinical models | 2013 | 10.1006/j.jmig.2013.01.018 | Gossion W | Stud Health Technol Inform | 2013 | |
| 2978769 | Reuse of public, genome-wide, mouse embryonic development data | Gracia JC | J Public Health | 2018 | 10.1007/s11524-018-0448-1 | 25848463 | Recommendations for the use of operational electronic health record data in comparative effectiveness research | 2020 | 10.1093/ibd/ibz009-1603104 | Alter G | Gastroenterology | 2020 | 10.1093/ibd/ibz009-1603104 |
| 3089574 | Impact of Electronic versus Paper-Based Recording before FHIR Implementation on Health Care Professionals' Perceptions of Data Quality, Data Quantity, and Data Reuse | Juskals E | Appl Clin Inform | 2019 | 10.1053/j.ajic.2019.10.004 | 25091808 | Stability metrics for multi-source biomedical data based on simplicial projections from probability distribution distances | 2017 | 10.1177/0962280214545122 | Sve'z C | Stat Methods Med Res | 2017 | 10.1177/0962280214545122 |
| 3030551 | Changes in Data Sharing and Data Access Practices and Perceptions among Scientists Worldwide | Thompson C | PLoS One | 2018 | 10.1371/journal.pone.0219420 | 24747879 | The need for harmonized structured documentation and chains of secondary use - results of a systematic analysis with automated form comparison for prostate and breast cancer | 2014 | 10.1016/j.jmig.2014.04.008 | Cosson S | Stud Health Technol Inform | 2019 | 10.1038/nrn.2016.08.004 |
| 3159447 | Compendium of cancer transformations for machine learning applications | Lim SB | So Data | 2019 | 10.1008/14597-019-0037-2 | 26152952 | A Query Tool Enabling Clinicians and Researchers to Explore Patient Cohorts | 2016 | 10.1038/nrn.2016.08.004 | Alter G | Gastroenterology | 2020 | 10.1093/ibd/ibz009-1603104 |
| 2895489 | Drug risk assessment and data reuse | Hyun S | Pharmacotherapy Drug Saf | 2018 | 10.1002/pds.318 | 28028995 | Evaluating the Paper-to-Screen Translation of Participant-Added Socioeconomic with High-Risk Participants | 2018 | 10.1186/s12913-018-03989-8 | Lim Chong Keung SN | Stud Health Technol Inform | 2015 | |
| 3032788 | Intelligence-based data warehouse environment to enable the reuse of electronic health record data | Canoy R | PLoS One | 2017 | 10.1371/journal.pone.0182426 | 29466022 | From Liguria HiV Web to Liguria Infectious Diseases Network: How a Digital Platform Empowered Doctors' Work and Patients' Care | 2018 | 10.1888/0944-2103-0109-5 | Hogon B | Proc SPIE Conf Hum Factor Comput Syst | 2016 | 10.1186/s12913-018-03989-8 |
| 3042079 | Kinematics of Big Biomedical Data to Characterize temporal variability and seasonality of data repositories: Functional Data Analysis of data temporal evolution over non-parametric statistical methods | Sve'z C | Int J Med Inform | 2018 | 10.1016/j.ijmedinf.2018.08.015 | 28625580 | Predicting biomedical metadata in CIDAR: A study of Gene Expression Omnibus (GEO) datasets | 2017 | 10.1038/nrn.2016.08.004 | Gianini B | ADIS Res Hum Retroviruses | 2016 | 10.1888/0944-2103-0109-5 |
| 29477819 | From Data Extraction to Analysis: Proposal of a Methodology to Optimize Health Data Reuse Process | Paulsson AD | Database (Oxford) | 2018 | 10.1093/database/bay020 | 28974262 | A generic method for improving the spatial interoperability of medical and ecological databases | 2016 | 10.1038/nrn.2016.08.004 | Chenias A | Int J Health Geogr | 2016 | 10.1038/nrn.2016.08.004 |
| 2732420 | Methodology and dimensionality of electronic health data quality assessment: enabling reuse for clinical research | Lamer A | Stud Health Technol Inform | 2018 | 10.1186/s12913-018-00861-1 | 28028995 | Evaluating the Paper-to-Screen Translation of Participant-Added Socioeconomic with High-Risk Participants | 2018 | 10.1186/s12913-018-03989-8 | Diaz-Garcia JF | AMIA J Summits Transl Soc Proc | 2015 | |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 25120292 | ProtonCNavigator: emulation based software for the design, documentation and reproduction of biological experiments | 2014 | 10.1186/1365-2754-11-113 | Chenias A | Int J Health Geogr | 2016 | 10.1038/nrn.2016.08.004 |
| 2900627 | Evidence-Based Guidelines for Intrafire Design for Data Entry in Electronic Health Records | Wilbanks BA | Comput Inform Nurs | 2018 | 10.1007/978-1-4939-9000-7_17 | 28974262 | A generic method for improving the spatial interoperability of medical and ecological databases | 2016 | 10.1038/nrn.2016.08.004 | Chenias A | Int J Health Geogr | 2016 | 10.1038/nrn.2016.08.004 |
| 3002178 | Intelligence-based data warehouse environment to enable the reuse of electronic health record data | Canoy R | PLoS One | 2017 | 10.1371/journal.pone.0182426 | 28028995 | Evaluating the Paper-to-Screen Translation of Participant-Added Socioeconomic with High-Risk Participants | 2018 | 10.1186/s12913-018-03989-8 | Diaz-Garcia JF | AMIA J Summits Transl Soc Proc | 2015 | |
| 3007796 | High Performance Monitor Estimation for Image Sensors with Video Compression | Xu W | Sensors (Basel) | 2015 | 10.3390/s15080792 | 28974262 | A generic method for improving the spatial interoperability of medical and ecological databases | 2016 | 10.1038/nrn.2016.08.004 | Chenias A | Int J Health Geogr | 2016 | 10.1038/nrn.2016.08.004 |
| 3024800 | Intelligence-based data warehouse environment to enable the reuse of electronic health record data | Canoy R | PLoS One | 2017 | 10.1371/journal.pone.0182426 | 28974262 | A generic method for improving the spatial interoperability of medical and ecological databases | 2016 | 10.1038/nrn.2016.08.004 | Chenias A | Int J Health Geogr | 2016 | 10.1038/nrn.2016.08.004 |
| 2719213 | MAUCBasis: A Tool to Support Data Reuse in OmicsResearch | Bonetto M | IEEE J Transl Eng Health Med | 2015 | 10.1109/JTEHM.2015.2510343 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2964470 | Controlled comparison of human transcriptional bioreactor data | Goldman NP | So Data | 2019 | 10.1008/14597-019-0037-2 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 2950120 | Three components in health care data reuse | Peres H | Nature | 2017 | 10.1038/5401586-0181-017314 | 24959874 | A patient-centered longitudinal care plan: vision vs. reality | 2014 | 10.1186/1365-2754-11-113 | Dykes PC | J Am Med Inform Assoc | 2014 | 10.1186/1365-2754-11-113 |
| 295 | | | | | | | | | | | | | |

Exemplos

- Utilização de diversos conjuntos de dados de amostras do sistema immune de ovinos
- Identificação de lncRNAs associados à resposta immune em ovinos



Exemplos



Referências

- Arend D, Lange M, Chen J, Colmsee C, Flemming S, Hecht D, Scholz U. 2014. e!DAL - a framework to store, share and publish research data. *BMC Bioinformatics*. 15(1):214.
- Bierer BE, Crosas M, Pierce HH. 2017. Data Authorship as an Incentive to Data Sharing. *New England Journal of Medicine*. 376(17):1684–1687.
- Bilbao-Arribas M, Jugo BM. 2022. Transcriptomic meta-analysis reveals unannotated long non-coding RNAs related to the immune response in sheep. *Front Genet*. 13. [accessed 2024 Jul 3]. <https://www.frontiersin.org/journals/genetics/articles/10.3389/fgene.2022.1067350/full>.
- Borghi J, Van Gulick A. 2022. Promoting Open Science Through Research Data Management. *Harvard Data Science Review*. [accessed 2024 Jul 3]. <https://hdsr.mitpress.mit.edu/pub/72kcw990>.
- Cunha-Oliveira T, Ioannidis JPA, Oliveira PJ. 2024. Best practices for data management and sharing in experimental biomedical research. *Physiological Reviews*. 104(3):1387–1408.
- Gomes DGE, Pottier P, Crystal-Ornelas R, Hudgins EJ, Foughirad V, Sánchez-Reyes LL, Turba R, Martinez PA, Moreau D, Bertram MG, et al. 2022. Why don't we share data and code? Perceived barriers and benefits to public archiving practices. *Proceedings of the Royal Society B: Biological Sciences*. 289(1987):20221113.
- Goodman A, Pepe A, Blocker AW, Borgman CL, Cranmer K, Crosas M, Stefano RD, Gil Y, Groth P, Hedstrom M, et al. 2014. Ten Simple Rules for the Care and Feeding of Scientific Data. *PLOS Computational Biology*. 10(4):e1003542.
- Hafner A, DeLeo V, Deng C, Elisk CG, Fleming D, Harrison PW, Kalbfleisch TS, Petry B, Pucker B, Quezada-Rodríguez EH, et al. 2024. Data Reuse in Agricultural Genomics Research: Present Challenges and Future Solutions. [accessed 2024 Jul 3]. <https://www.preprints.org/manuscript/202401.0780/v1>.
- McDonald-McGinn DM, Sullivan KE. 2011. Chromosome 22q11.2 Deletion Syndrome (DiGeorge Syndrome/Velocardiofacial Syndrome). *Medicine*. 90(1):1.
- Mello MM, Lieou V, Goodman SN. 2018. Clinical Trial Participants' Views of the Risks and Benefits of Data Sharing. *New England Journal of Medicine*. 378(23):2202–2211.
- Raman A. 2021. Secondary Data Analysis: Lessons and perspective of a research parasite. [accessed 2024 Jul 23]. <https://osf.io/bec69>.
- Schwab S, Janiaud P, Dayan M, Amrhein V, Panczak R, Palagi PM, Hemkens LG, Ramon M, Rothen N, Senn S, et al. 2022. Ten simple rules for good research practice. *PLOS Computational Biology*. 18(6):e1010139.
- Sielemann K, Hafner A, Pucker B. 2020. The reuse of public datasets in the life sciences: potential risks and rewards. *PeerJ*. 8:e9954.
- Yuan H, Hicks P, Ahmadian M, Johnson K, Valtadoros L, Krishnan A. 2024. Annotating publicly-available samples and studies using interpretable modeling of unstructured metadata. [accessed 2024 Jul 3]. <http://biorxiv.org/lookup/doi/10.1101/2024.06.03.597206>.

Obrigado!

carvalhopc@usp.br
linktr.ee/carvalhopc