

GOOD PRACTICE IN HUMAN DATA HANDLING AND ANALYSIS

IMFORFUTURE MEETING 2018, ZAGREB

Lennart Karssen

26–28 March 2018

PolyΩmica

's-Hertogenbosch, NL

l.c.karssen@polyomica.com

1. Introduction
2. Data stewardship
 - Data stewardship, what, why, how
 - Infrastructure
 - Provenance
 - FAIR principles
 - Reproducible research
3. Ethics
4. Conclusions/wrapping up

INTRODUCTION

ABOUT ME

- Dutch
- 2008: PhD in experimental atomic physics, Utrecht University, NL
- 2009: Teaching physics @ National University of Rwanda
- 2009: Unix/Linux consultant @ Snow B.V., NL
- 2010: Postdoc bioinformatics, data manager, sysadmin @ genetic epidemiology group, Erasmus MC, Rotterdam, NL
- 2013: PolyOmica
- 2018: PolyKnomics

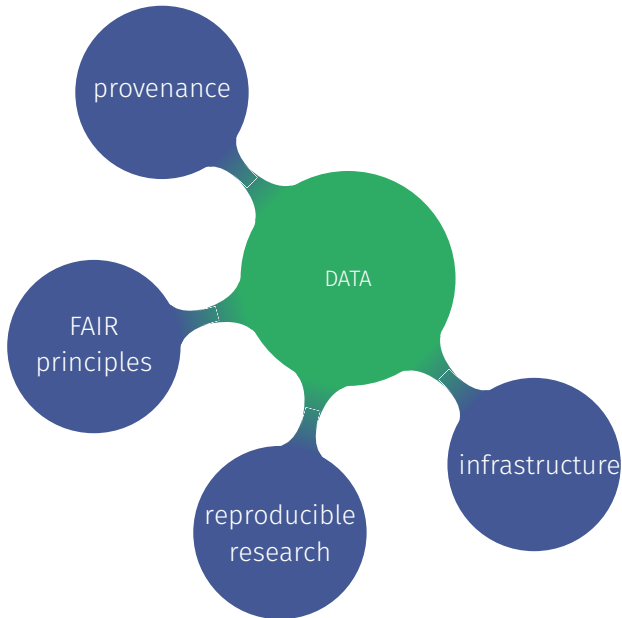


Image source: Wikipedia, snow.nl, erasmusmc.nl

- Research
 - omics analysis: genomics, glycomics,
 - development of statistical methods
 - scientific software development
 - EU-funded projects like IMforFUTURE, PainOmics
- (Scientific) consulting
 - experimental design
 - analyses
 - compute/storage infrastructure
- Education
 - (statistical) genetics
 - epidemiology
 - scientific computing
 - ...

DATA STEWARDSHIP

WHAT IS DATA STEWARDSHIP?



WHAT IS DATA STEWARDSHIP?

- Lots of data generated: in the lab, on computers: Big Data/Data Tsunami

WHAT IS DATA STEWARDSHIP?

- Lots of data generated: in the lab, on computers: Big Data/Data Tsunami
- What about
 - the quality of the data?
 - analysis?
 - the reproducibility of research?
 - re-use of new and existing data?

WHAT IS DATA STEWARDSHIP?

- Lots of data generated: in the lab, on computers: Big Data/Data Tsunami
- What about
 - the quality of the data?
 - analysis?
 - the reproducibility of research?
 - re-use of new and existing data?

In short

The past, present and future quality and availability of data

- Where to store the data for analysis?
- Where to archive?
- What to archive?
- Where to run analyses?
- Who pays?
- Shouldn't forget security/privacy!

Aim for

- proper collection
- annotation
- archiving
- long-term care

of data.

Aim for

- proper collection
- annotation
- archiving
- long-term care

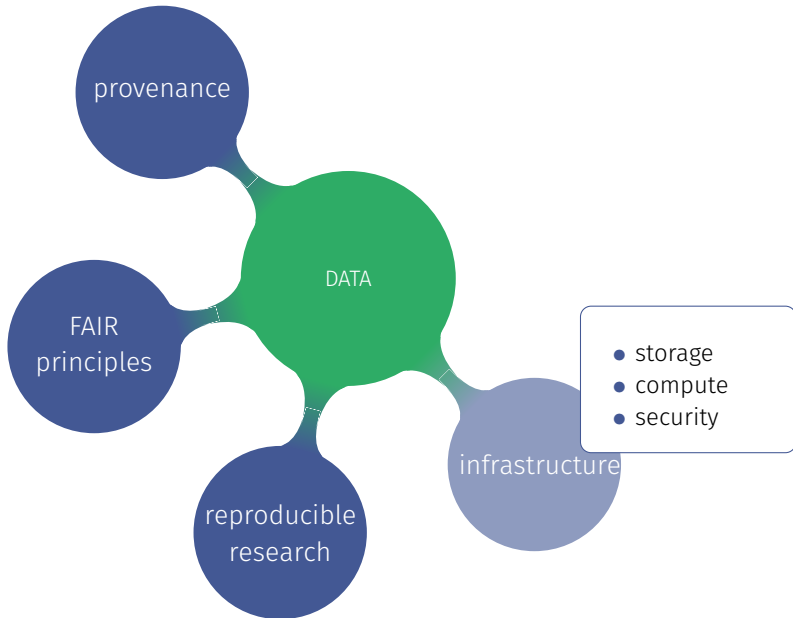
of data. Not just traditional data, also

- algorithms
- tools
- workflows

leading to the data

- Funding agencies require data management plan
- Funding agencies require integration with existing data
- Publishers require access to (some of) the data

WHAT IS DATA STEWARDSHIP?



SOME NUMBERS

Disk space usage of common data sets

Data set	Size
<i>Lipidomics GWAS: 117 lipids, 20M SNPs, 3.5k people</i>	122 GB
<i>1kG imputed data: 30M SNPs, 3.5k people</i>	500 GB
<i>UK Biobank imputed data: 90M SNPs, 500k people</i>	2.1TB
<i>Aligned NGS data (BAM files): 1.3k people</i>	7.2 TB

MORE NUMBERS

Computation time

Data set	Time
<i>Lipidomics GWAS: 117 lipids, 20M SNPs, 3.5k people</i>	10 d
<i>1kG imputation: 30M SNPs, 3.5k people</i>	1 w
<i>Aligning NGS data (creating BAM files): 1.3k people</i>	weeks

MORE NUMBERS

Computation time

Data set	Time
<i>Lipidomics GWAS: 117 lipids, 20M SNPs, 3.5k people</i>	10 d
<i>1kG imputation: 30M SNPs, 3.5k people</i>	1 w
<i>Aligning NGS data (creating BAM files): 1.3k people</i>	weeks

Missing information:

How many CPUs/compute nodes were used?

MORE NUMBERS

Computation time

Data set	Time
<i>Lipidomics GWAS: 117 lipids, 20M SNPs, 3.5k people</i>	10 d
<i>1kG imputation: 30M SNPs, 3.5k people</i>	1 w
<i>Aligning NGS data (creating BAM files): 1.3k people</i>	weeks

Missing information:

How many CPUs/compute nodes were used?

Transferring data over a network

1TB via 1 Gbps (gigabit per second) network: ≈ 2.2 h

MORE NUMBERS

Computation time

Data set	Time
<i>Lipidomics GWAS: 117 lipids, 20M SNPs, 3.5k people</i>	10 d
<i>1kG imputation: 30M SNPs, 3.5k people</i>	1 w
<i>Aligning NGS data (creating BAM files): 1.3k people</i>	weeks

Missing information:

How many CPUs/compute nodes were used?

Transferring data over a network

1TB via 1 Gbps (gigabit per second) network: ≈ 2.2 h

theoretically...

MORE NUMBERS: COSTS OF SERVER-CLASS HARDWARE

The "mini cluster" of one of our clients

Year	Item	Price (€)
2009	new compute + storage server	30k
2010	add compute + storage server	21k
2013	add compute nodes + extra storage	34k
2016	replace compute + storage server	11k
2017	replace/increase storage	11k
2018	(?) increase storage + compute	~25k

MORE NUMBERS: COSTS OF SERVER-CLASS HARDWARE

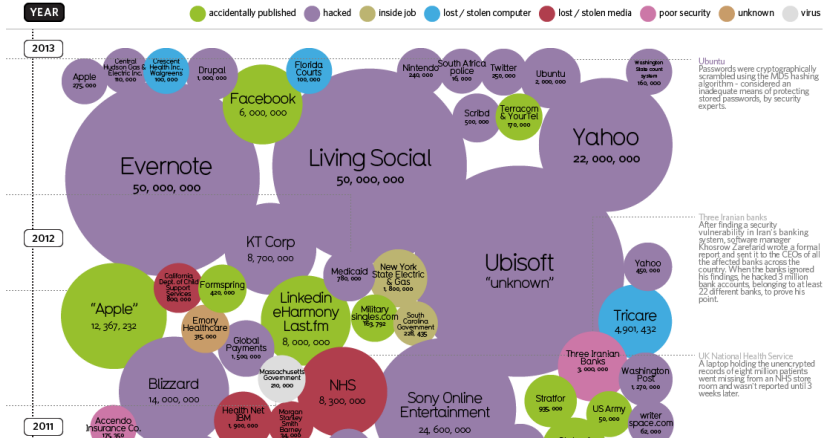
The "mini cluster" of one of our clients

Year	Item	Price (€)
2009	new compute + storage server	30k
2010	add compute + storage server	21k
2013	add compute nodes + extra storage	34k
2016	replace compute + storage server	11k
2017	replace/increase storage	11k
2018	(?) increase storage + compute	~25k

Excluding

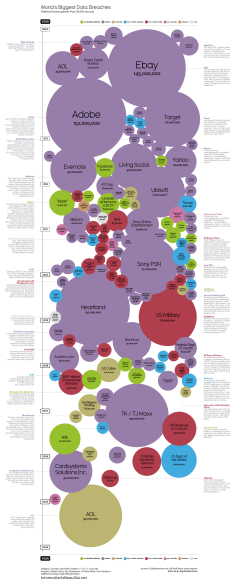
- installation costs
- maintenance costs
- data centre costs

SECURITY: DATA BREACHES



© Information is Beautiful.net

SECURITY: DATA BREACHES



Revealed: 50 million Facebook profiles harvested for Cambridge Analytica in major data breach

How Trump Consultants Exploited the Facebook Data of Millions

**A HURRICANE FLATTENS
FACEBOOK**

- Secure connections: SSH, HTTPS
- Don't share accounts or passwords
- Use different, long, complicated passwords for each server/service/site
 - → use a password manager (e.g. KeePass, KeePassX, LastPass, ...)
- Use SSH keys for logging in to remote servers
- Keep names, addresses, post codes, etc. and the link to sample IDs separate in a safe location

- Secure connections: SSH, HTTPS
- Don't share accounts or passwords
- Use different, long, complicated passwords for each server/service/site
 - → use a password manager (e.g. KeePass, KeePassX, LastPass, ...)
- Use SSH keys for logging in to remote servers
- Keep names, addresses, post codes, etc. and the link to sample IDs separate in a safe location

Actually...

Unless you are the data manager: hand over to the data manager and delete it!

SECURITY: TIPS (CONTD.)

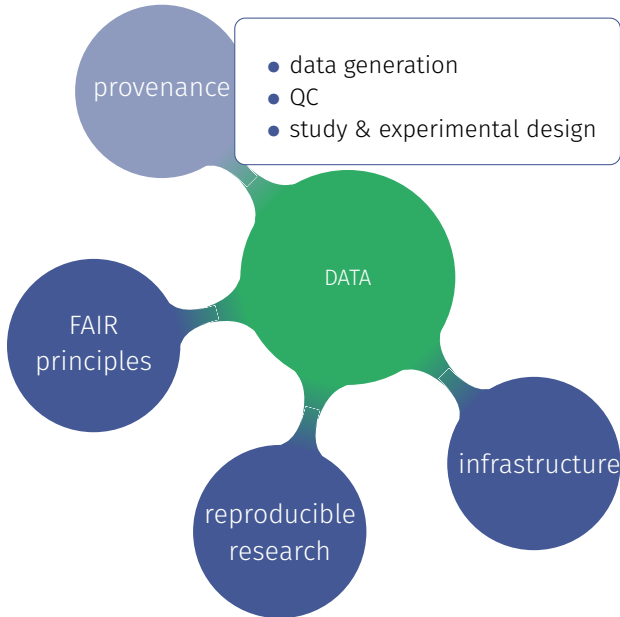
- Store data only on servers or PCs at work

- Store data only on servers or PCs at work

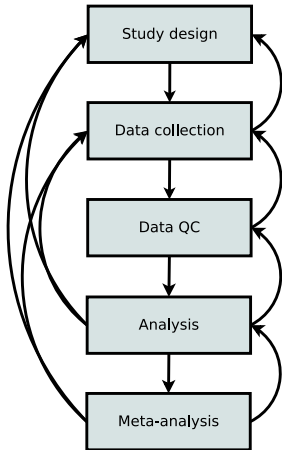
Please, not on...

- USB sticks
- External harddisks
- Laptops
- DropBox/MS OneDrive/Google Drive
 - You loose control:
 - delete != gone
 - you may give them the right to use and/or modify the data
 - service may shut down
 - "There is not such thing as the cloud. Only somebody else's computer"
- (unless sufficiently encrypted)

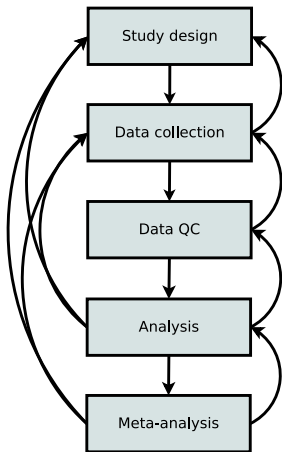
WHAT IS DATA STEWARDSHIP?



A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



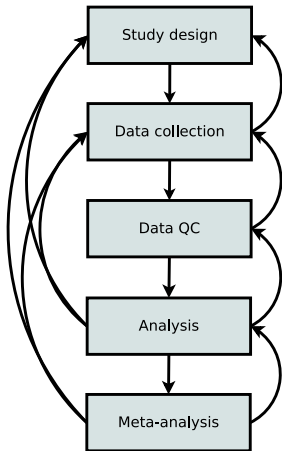
A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



Study design

- What data to collect? Where? From whom? For how long?
- Staff

A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



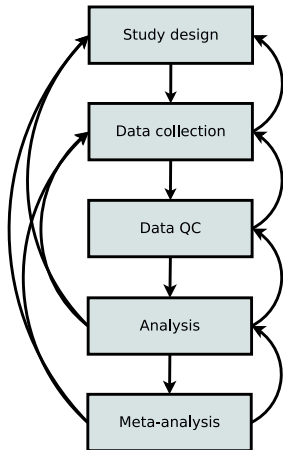
Study design

- What data to collect? Where? From whom? For how long?
- Staff

Data collection

- Plasma, MRI, blood pressure, questionnaires, ...
- Medical doctors, biologists, technicians, data manager, *you*?

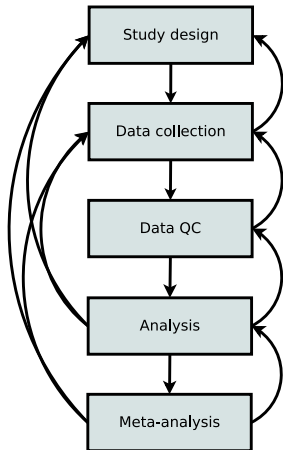
A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



Quality Control

- Remove outliers and other crazy values, check validity
- Units of measured values (!)
- Requires expert knowledge of the data
- Post-docs, data manager, *you?*

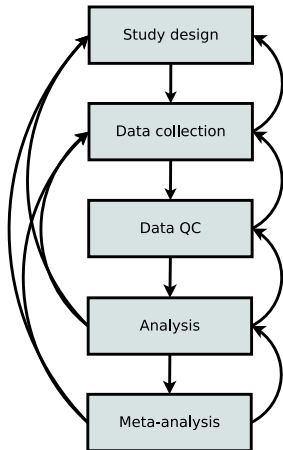
A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



Analysis

- use the data for research
- post-docs, PhD students, staff, *you*

A TYPICAL WORKFLOW IN "OMICS EPIDEMIOLOGY"



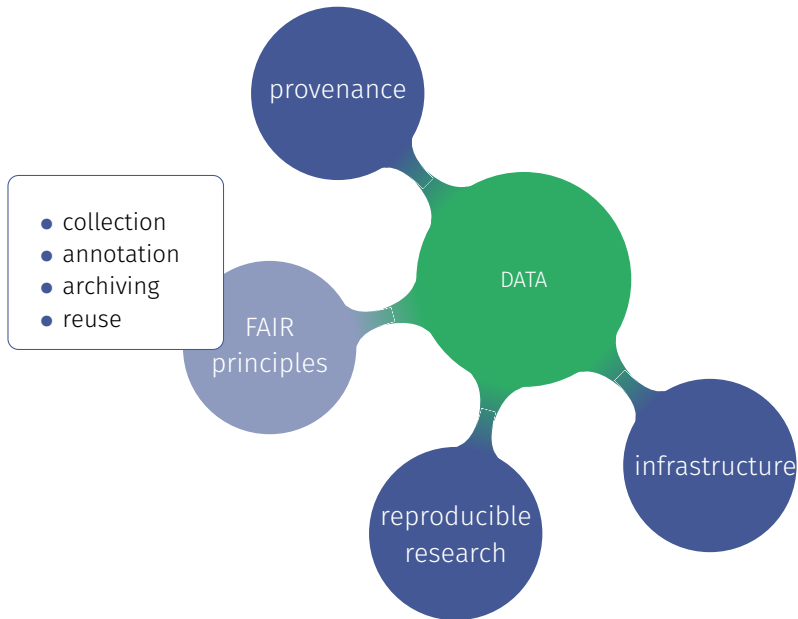
Analysis

- use the data for research
- post-docs, PhD students, staff, *you*

Meta-analysis

- combine data from various cohorts; harmonization
- post-docs, PhD students, *you*

WHAT IS DATA STEWARDSHIP?



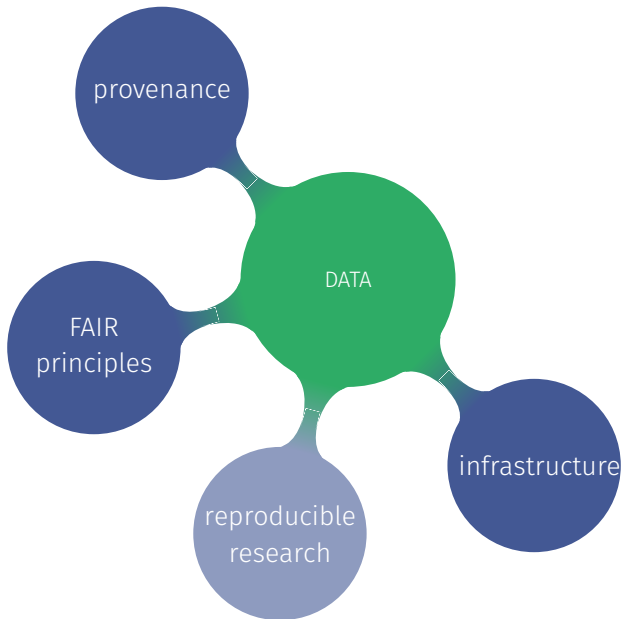
- F - Findable
- A - Accessible
- I - Interoperable
- R - Reusable

doi:10.1038/sdata.2016.18

“Good data management is **not a goal in itself**, but rather is the key conduit leading to **knowledge discovery** and **innovation**, and to subsequent **data and knowledge integration** and **reuse** by the community after the data publication process”

- emphasis on machine readability
- rich meta data
- open protocols/file formats
- links/references
- license

WHAT IS DATA STEWARDSHIP?



- Transparency (recent claims of data manipulation!)
- Whatever you do should be reproducible for
 - yourself
 - rerun with different parameters
 - after X months when writing a paper/thesis
 - your successor(s)
 - your supervisors
 - journals/your peers

- Transparency (recent claims of data manipulation!)
- Whatever you do should be reproducible for
 - yourself
 - rerun with different parameters
 - after X months when writing a paper/thesis
 - your successor(s)
 - your supervisors
 - journals/your peers

How?

- Documentation
- Version control
- Literate programming

ETHICS

(RECENT) HISTORY

- 1939–1945: Nazi medical war crimes
 - Outcome: Nuremberg Code → First international code of research ethics
- 1948: UN Universal Declaration of Human Rights



Image source: Wikipedia

1930–1972: Syphilis study at Tuskegee (Alabama, USA)

- ~ 600 African-American men (~ 400 cases, ~ 200 controls)
- No informed consent
- some procedures (e.g. spinal taps): "special free treatment"
- continued even after treatment (penicillin) was found in the 1940s
- Outcomes: national policies, medical-ethical committees

1944–1974: Cold war human radiation experiments

- > 400 experiments to determine the effects of radiation on human health; calibration of instruments for radiation detection
- Most studies involved minimal risks, appropriate informed consent
- However, some subject suffered injuries in studies that offered no prospect of direct benefit, or were presented as standard practice

Three principles essential to ethical conduct of research with humans

- Respect for persons
 - individuals should be treated as autonomous agents
 - persons with diminished autonomy are entitled to additional protections
- Beneficence
 - do no harm
 - maximise possible benefits, minimise possible harms
- Justice
 - individuals and groups should be treated fairly and equitably w.r.t. burdens and benefits of the research

- Consider:
 - vulnerable populations needing specific protection (e.g. children, prisoners)
 - risks associated in participation in research and appropriate protections
 - inclusion/exclusion criteria
- Informed consent
- Medical-ethical review board

Once we have blood/DNA...

- Broad consent?
- Explicitly ask permission for each research question?
- How about unforeseen questions?
- How about DNA vs. privacy?
- Opt out?

OPEN QUESTIONS FOR OUR FIELD

Once we have blood/DNA...

- Broad consent?
- Explicitly ask permission for each research question?
- How about unforeseen questions?
- How about DNA vs. privacy?
- Opt out?

Open data

- Should clinical trial data be open?
- Should biobank data be open?
- (and what does 'open' mean?)

EU GENERAL DATA PROTECTION REGULATION (GDPR)

- enforceable from 25 May 2018
- single set of rules for all EU member states
- if consent is used as basis for processing:
 - must be explicit for data collected, and the purposes data is used for
 - consent for children must be given by the child's parent or custodian, and verifiable
 - data controllers must be able to prove "consent" (opt-in)
 - consent may be withdrawn
- data breach: data controller has to notify the authorities without undue delay
- right to erasure
- data portability: transfer personal data from one electronic processing system to and into another
- much more...

CONCLUSIONS/WRAPPING UP

WHY SHOULD YOU CARE ABOUT DATA STEWARDSHIP?

Privacy

- Personal data from real people
- Genetic data: also affects their offspring
- Moreover, data leaks:
 - are bad publicity
 - may affect other people's willingness to participate
 - ...

WHY SHOULD YOU CARE ABOUT DATA STEWARDSHIP?

Privacy

- Personal data from real people
- Genetic data: also affects their offspring
- Moreover, data leaks:
 - are bad publicity
 - may affect other people's willingness to participate
 - ...

Cost of creating and maintaining a data set

- Sample collection (interviews, biosamples, ...)
- Curation
- Storage (not just digital, but freezers, cell lines, ...)

High-quality research

- Understand your data
 - What was measured?
 - How was it measured?
 - What are reasonable values? Consult experts! Use units!
 - Understand data formats (e.g. VCF vs. plink)
- Understand your software: what does it do?
- Provenance
 - Where did the data come from?
 - Authoritative source?
 - What QC was done?

- Take good care of the data you work with (now and in the future)
- Understand your data: provenance, ethics
- Document what you do
- Talk to people → integration
- When you apply for grants, don't forget budget for:
 - data manager
 - bioinformatics/IT infrastructure

THANK YOU!

Any questions?

THANK YOU!

Any questions?