

«Human-Centered Data Science»
**Enabling Reproducibility
of Your Data Science Practice**

Prof. Dr. Claudia Müller-Birn
Human-Centered Computing, Institute of Computer Science
Freie Universität Berlin
May 12, 2022

Lecture Overview

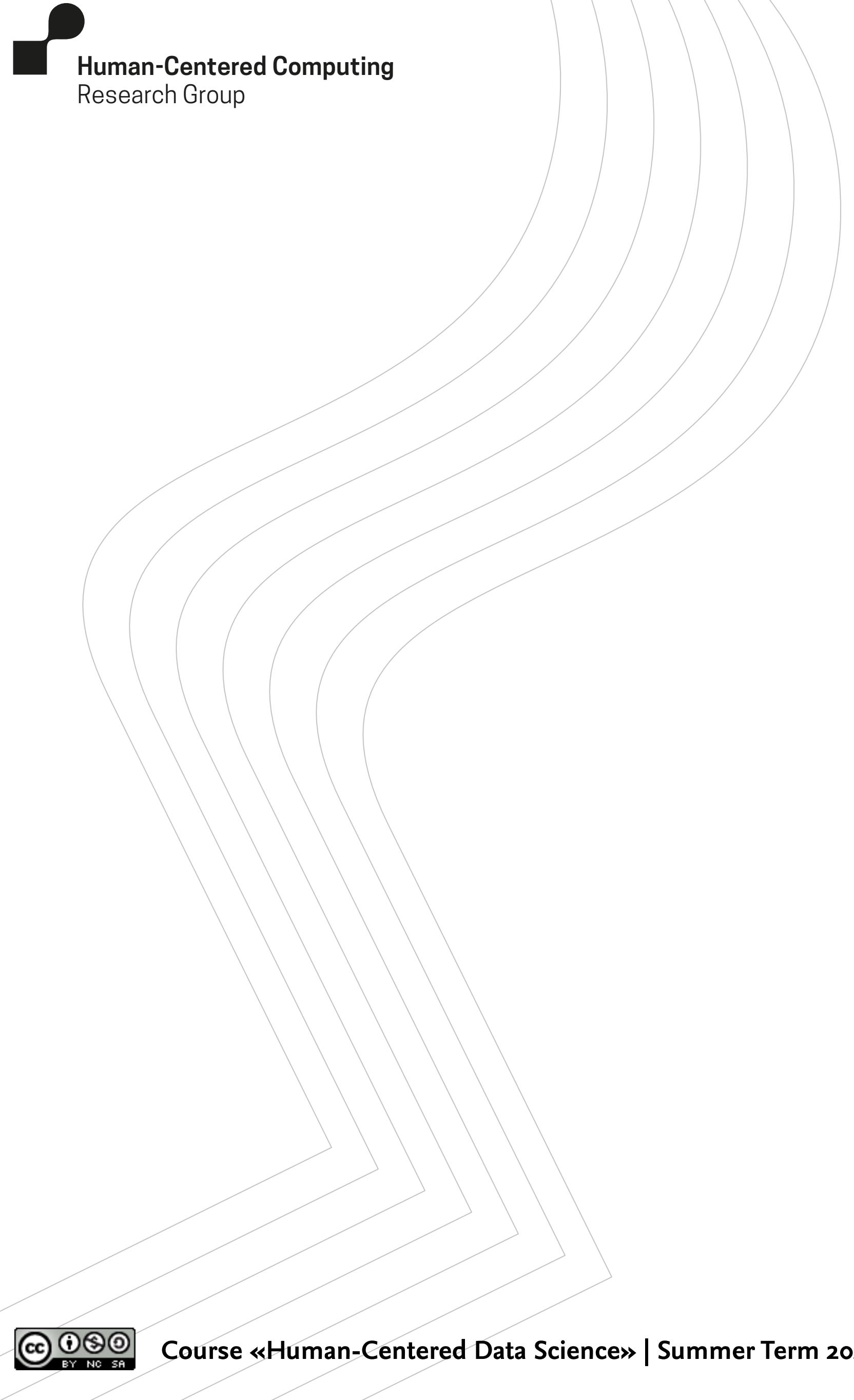
Recap

Theoretical Considerations of Accountability (aspects of accountability, relationships, professional accountability = reproducibility, definition of repeatability, reproducibility and replicability, degree of replication)

☕ Break

Translating into Practice (recommendations for sharing, using, and understanding your data science workflows, stages of a basic reproducible workflow and things to think of)





Recap



What is Ethics?

“

“Ethics as a discipline is the philosophical and systematic study of morality which is expressed in enquiries into the nature of the good life, how we should live, what kind of society we want to live in, and how we should treat others.”

Ethics is a philosophical discipline.

Morality is a set of moral norms, feelings, attitudes, or actions.

The Global Landscape of AI Ethics Guidelines

Table 1 | Ethics guidelines for AI by country of issuer (Australia–UK)

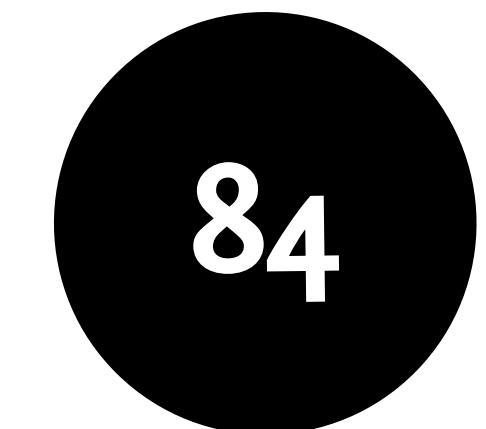
Name of document/website	Issuer	Country of issuer
Artificial Intelligence: Australia's Ethics Framework: A Discussion Paper	Department of Industry Innovation and Science	Australia
Montréal Declaration: Responsible AI	Université de Montréal	Canada
Work in the Age of Artificial Intelligence: Four Perspectives on the Economy, Employment, Skills and Ethics	Ministry of Economic Affairs and Employment	Finland
Tieto's AI Ethics Guidelines	Tieto	Finland
Commitments and Principles	OP Group	Finland
How Can Humans Keep the Upper Hand? Report on the Ethical Matters Raised by AI Algorithms	French Data Protection Authority (CNIL)	France
For a Meaningful Artificial Intelligence: Towards a French and European Strategy	Mission Villani	France
Ethique de la Recherche en Robotique	CERNA (Allistene)	France
AI Guidelines	Deutsche Telekom	Germany
SAP's Guiding Principles for Artificial Intelligence	SAP	Germany
Automated and Connected Driving: Report	Federal Ministry of Transport and Digital Infrastructure, Ethics Commission	Germany
Ethics Policy	Icelandic Institute for Intelligent Machines (IIIM)	Iceland
Discussion Paper: National Strategy for Artificial Intelligence	National Institution for Transforming India (NITI Aayog)	India
L'intelligenza Artificiale al Servizio del Cittadino	Agenzia per l'Italia Digitale (AGID)	Italy
The Japanese Society for Artificial Intelligence Ethical Guidelines	Japanese Society for Artificial Intelligence	Japan
Report on Artificial Intelligence and Human Society (unofficial translation)	Advisory Board on Artificial Intelligence and Human Society (initiative of the Minister of State for Science and Technology Policy)	Japan
Draft AI R&D Guidelines for International Discussions	Institute for Information and Communications Policy (IICP), The Conference toward AI Network Society	Japan
Sony Group AI Ethics Guidelines	Sony	Japan
Human Rights in the Robot Age Report	The Rathenau Institute	Netherlands
Dutch Artificial Intelligence Manifesto	Special Interest Group on Artificial Intelligence (SIGAI), ICT Platform Netherlands (IPN)	Netherlands
Artificial Intelligence and Privacy	The Norwegian Data Protection Authority	Norway
Discussion Paper on Artificial Intelligence (AI) and Personal Data—Fostering Responsible Development and Adoption of AI	Personal Data Protection Commission Singapore	Singapore
Mid- to Long-Term Master Plan in Preparation for the Intelligent Information Society	Government of the Republic of Korea	South Korea
AI Principles of Telefónica	Telefónica	Spain
AI Principles & Ethics	Smart Dubai	UAE
Principles of robotics	Engineering and Physical Sciences Research Council UK (EPSRC)	UK
The Ethics of Code: Developing AI for Business with Five Core Principles	Sage	UK
Big Data, Artificial Intelligence, Machine Learning and Data Protection	Information Commissioner's Office	UK
DeepMind Ethics & Society Principles	DeepMind Ethics & Society	UK
Business Ethics and Artificial Intelligence	Institute of Business Ethics	UK
AI in the UK: Ready, Willing and Able?	UK House of Lords, Select Committee on Artificial Intelligence	UK
Artificial Intelligence (AI) in Health	Royal College of Physicians	UK
Initial Code of Conduct for Data-Driven Health and Care Technology	UK Department of Health & Social Care	UK
Ethics Framework: Responsible AI	Machine Intelligence Garage Ethics Committee	UK
The Responsible AI Framework	PricewaterhouseCoopers UK	UK
Responsible AI and Robotics: An Ethical Framework	Accenture UK	UK
Machine Learning: The Power and Promise of Computers that Learn by Example	The Royal Society	UK
Ethical, Social, and Political Challenges of Artificial Intelligence in Health	Future Advocacy	UK

Table 2 | Ethics guidelines for AI by country of issuer (USA, international, EU and N/A)

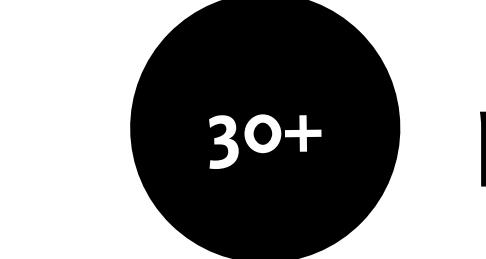
Name of document/website	Issuer	Country of issuer
Unified Ethical Frame for Big Data Analysis. IAF Big Data Ethics Initiative, Part A	The Information Accountability Foundation	USA
The AI Now Report: The Social and Economic Implications of Artificial Intelligence Technologies in the Near-Term	AI Now Institute	USA
Statement on Algorithmic Transparency and Accountability	Association for Computing Machinery (ACM)	USA
AI Principles	Future of Life Institute	USA
AI—Our Approach	Microsoft	USA
Artificial Intelligence: The Public Policy Opportunity	Intel Corporation	USA
IBM's Principles for Trust and Transparency	IBM	USA
OpenAI Charter	OpenAI	USA
Our Principles	Google	USA
Policy Recommendations on Augmented Intelligence in Health Care H-480.940	American Medical Association (AMA)	USA
Everyday Ethics for Artificial Intelligence: A Practical Guide for Designers and Developers	IBM	USA
Governing Artificial Intelligence: Upholding Human Rights & Dignity	Data & Society	USA
Intel's AI Privacy Policy White Paper: Protecting Individuals' Privacy and Data in the Artificial Intelligence World	Intel Corporation	USA
Introducing Unity's Guiding Principles for Ethical AI—Unity Blog	Unity Technologies	USA
Digital Decisions	Center for Democracy & Technology	USA
Science, Law and Society (SLS) Initiative	The Future Society	USA
AI Now 2018 Report	AI Now Institute	USA
Responsible Bots: 10 Guidelines for Developers of Conversational AI	Microsoft	USA
Preparing for the Future of Artificial Intelligence	Executive Office of the President; National Science and Technology Council; Committee on Technology	USA
The National Artificial Intelligence Research and Development Strategic Plan	National Science and Technology Council; Networking and Information Technology Research and Development Subcommittee	USA
AI Now 2017 Report	AI Now Institute	USA
Position on Robotics and Artificial Intelligence	The Greens (Green Working Group Robots)	EU
Report with Recommendations to the Commission on Civil Law Rules on Robotics	European Parliament	EU
Ethics Guidelines for Trustworthy AI	High-Level Expert Group on Artificial Intelligence	EU
AI4People—An Ethical Framework for a Good AI Society: Opportunities, Risks, Principles, and Recommendations	AI4People	EU
European Ethical Charter on the Use of Artificial Intelligence in Judicial Systems and Their Environment	Council of Europe: European Commission for the Efficiency of Justice (CEPEJ)	EU
Statement on Artificial Intelligence, Robotics and 'Autonomous' Systems	European Commission, European Group on Ethics in Science and New Technologies	EU
Artificial Intelligence and Machine Learning: Policy Paper	Internet Society	International
Report of COMEST on Robotics Ethics	COMEST/UNESCO	International
Ethical Principles for Artificial Intelligence and Data Analytics	Software & Information Industry Association (SIIA), Public Policy Division	International
ITI AI Policy Principles	Information Technology Industry Council (ITI)	International
Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems, Version 2	Institute of Electrical and Electronics Engineers (IEEE), The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems	International
Top 10 Principles for Ethical Artificial Intelligence	UNI Global Union	International
The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation	Future of Humanity Institute; University of Oxford; Centre for the Study of Existential Risk; University of Cambridge; Center for a New American Security; Electronic Frontier Foundation; OpenAI	International
White Paper: How to Prevent Discriminatory Outcomes in Machine Learning	WEF, Global Future Council on Human Rights 2016–2018	International

Table 2 | Ethics guidelines for AI by country of issuer (USA, international, EU and N/A) (Continued)

Name of document/website	Issuer	Country of issuer
Privacy and Freedom of Expression in the Age of Artificial Intelligence	Privacy International & Article 19	International
The Toronto Declaration: Protecting the Right to Equality and Non-discrimination in Machine Learning Systems	Access Now; Amnesty International	International
Charlevoix Common Vision for the Future of Artificial Intelligence	Leaders of the G7	International
Artificial Intelligence: Open Questions About Gender Inclusion	W20	International
Declaration on Ethics and Data Protection in Artificial Intelligence	ICDPPC	International
Universal Guidelines for Artificial Intelligence	The Public Voice	International
Ethics of AI in Radiology: European and North American Multisociety Statement	American College of Radiology; European Society of Radiology; Radiology Society of North America; Society for Imaging Informatics in Medicine; European Society of Medical Imaging Informatics; Canadian Association of Radiologists; American Association of Physicists in Medicine	International
Ethically Aligned Design: A Vision for Prioritizing Human Well-being with Autonomous and Intelligent Systems, First Edition (EAD1e)	Institute of Electrical and Electronics Engineers (IEEE), The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems	International
Tenets	Partnership on AI	N/A
Principles for Accountable Algorithms and a Social Impact Statement for Algorithms	Fairness, Accountability, and Transparency in Machine Learning (FATML)	N/A
10 Principles of Responsible AI	Women Leading in AI	N/A



guidelines



principles

Jobin, Anna, Marcello lenca, und Effy Vayena. „The global landscape of AI ethics guidelines“. Nature Machine Intelligence 1, Nr. 9 (1. September 2019): 389–99. <https://doi.org/10.1038/s42256-019-0088-2>.





Overview on Principles

Professional Responsibility

Fairness and Non-discrimination

Safety and Security

Privacy

CATEGORIES OF AI PRINCIPLES

Human Rights

Promotion of Human Values

Professional Responsibility

Human Control of Technology

Fairness and Non-discrimination

Transparency and Explainability

Safety and Security

Accountability

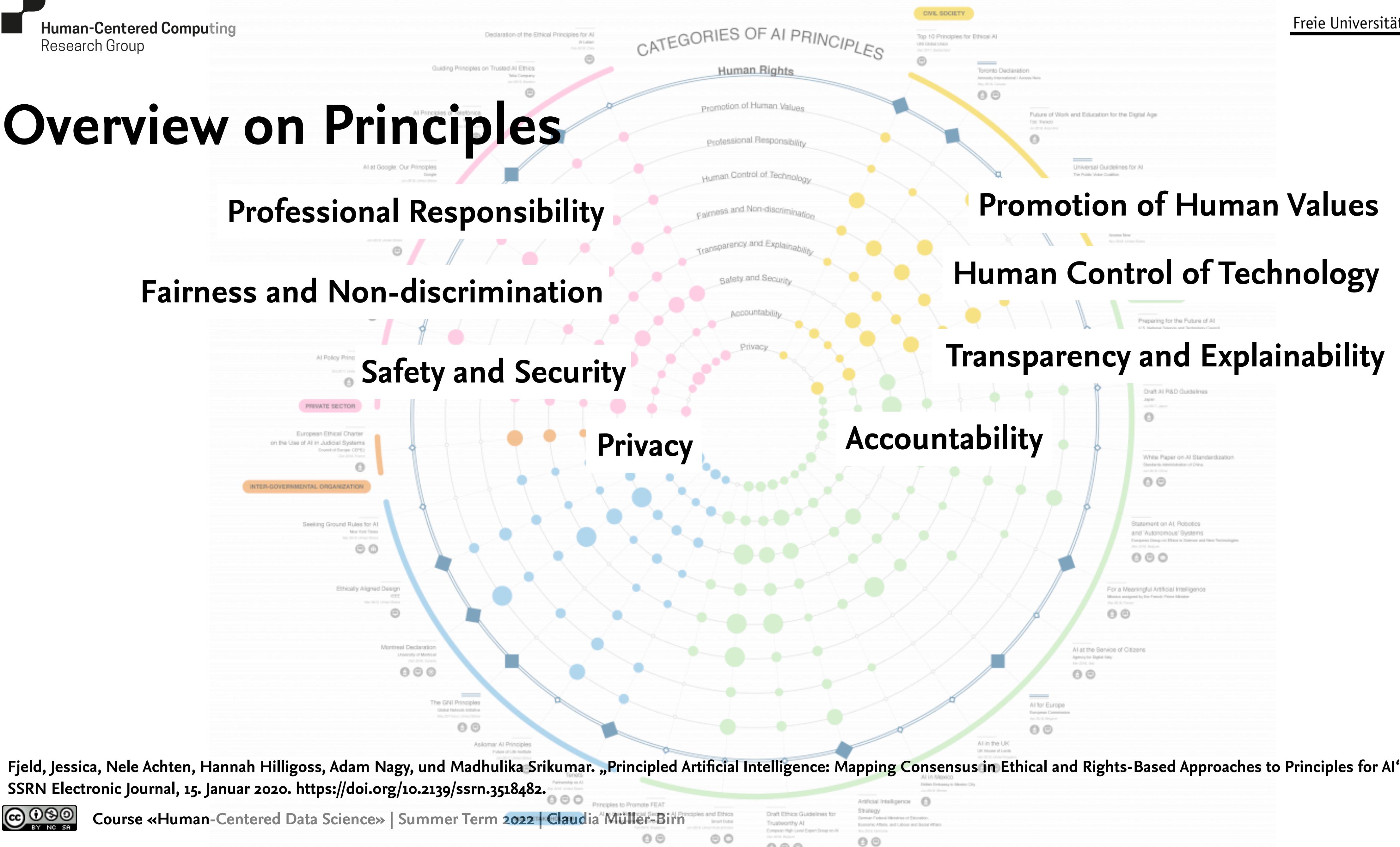
Privacy

Accountability

Promotion of Human Values

Human Control of Technology

Transparency and Explainability



Incorporating Principles in Your Data Science Practice

Principles

Professional Responsibility

Promotion of Human Values

Fairness and Non-discrimination

Human Control of Technology

Safety and Security

Transparency and Explainability

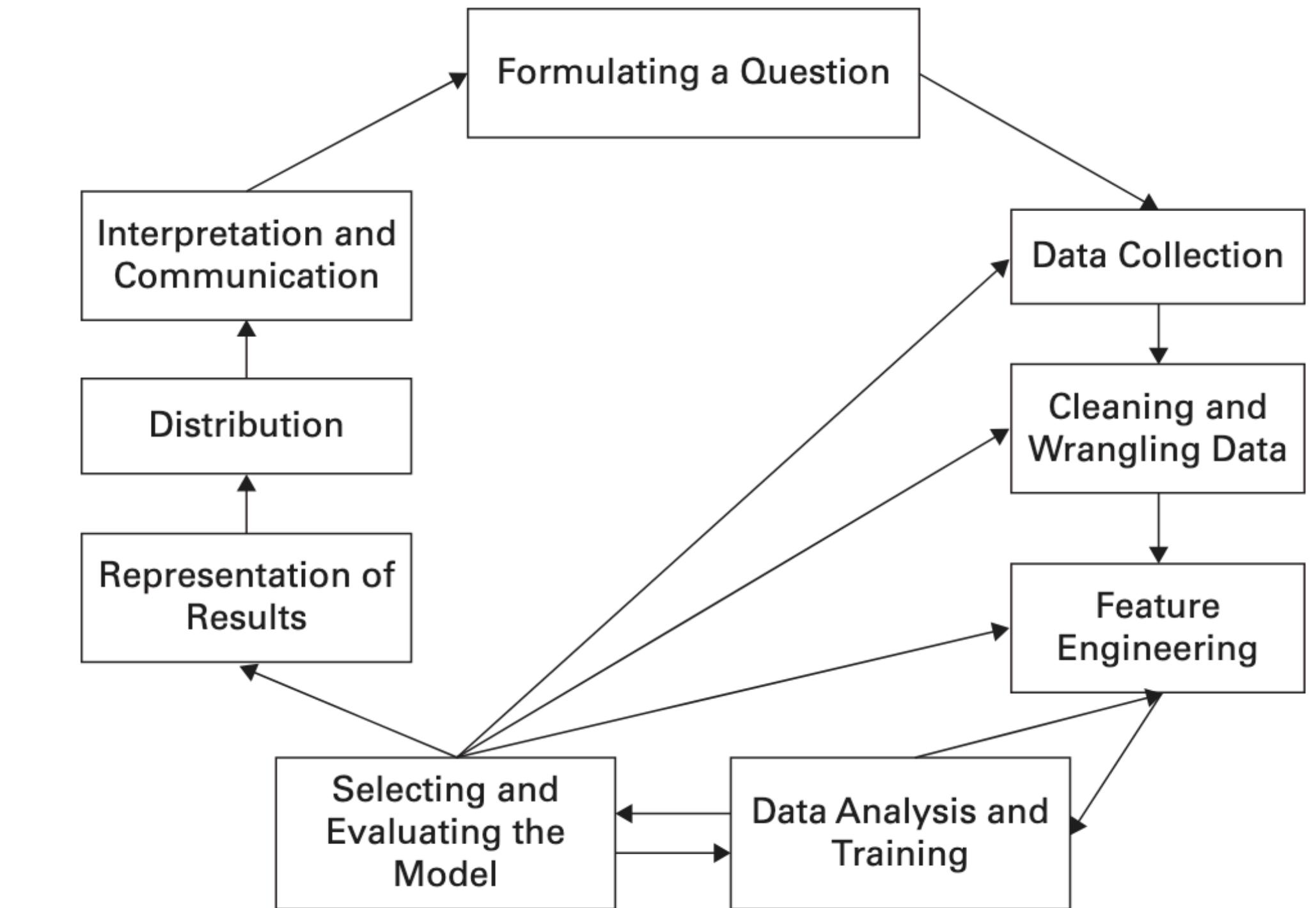
Accountability

Privacy

Values

Development of a shared ethos

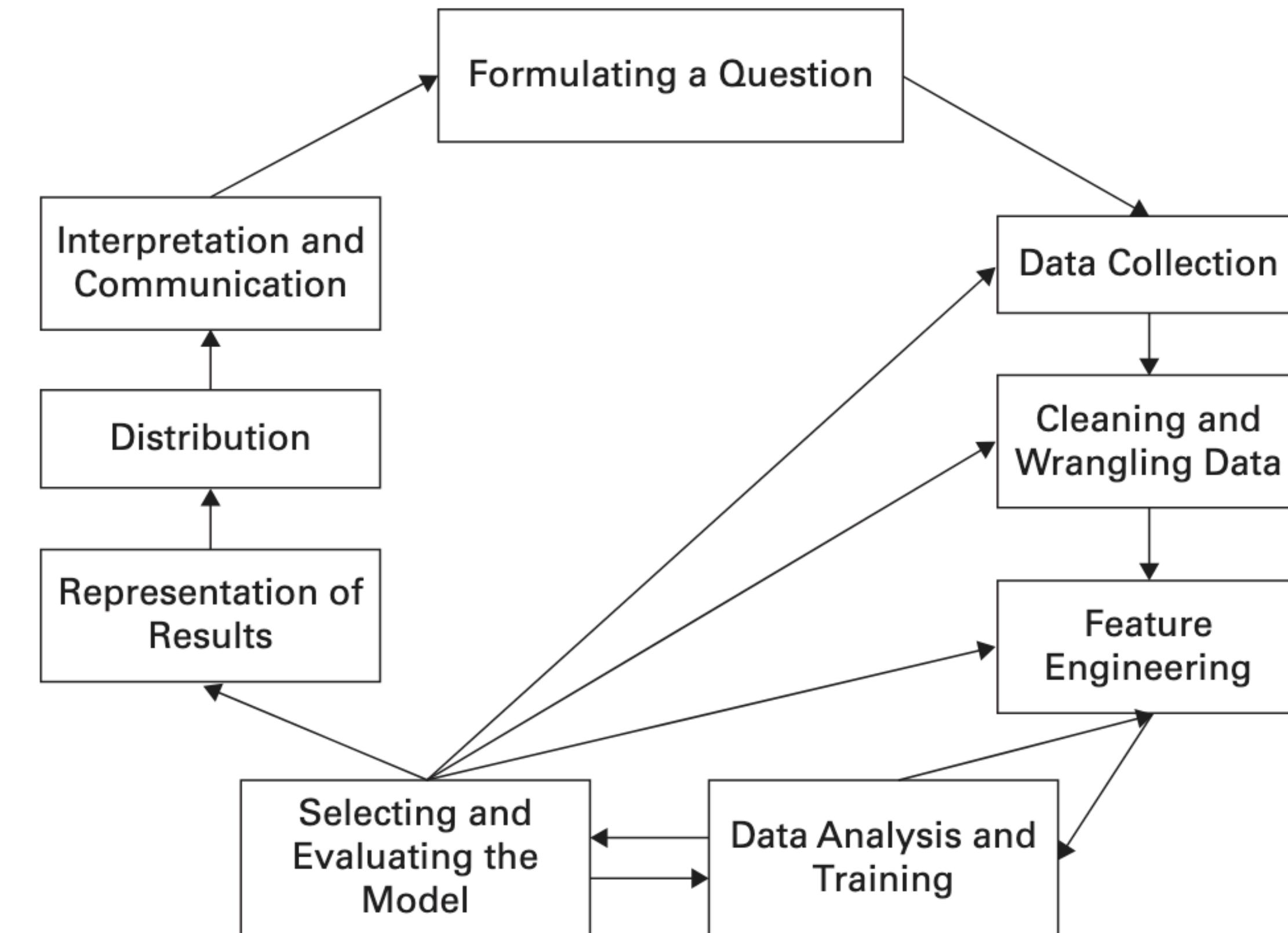
Values Levers

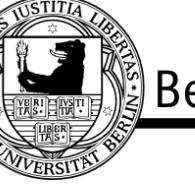


Incorporating Principles in Your Data Science Practice

Principles

- Professional Responsibility
- Promotion of Human Values
- Fairness and Non-discrimination
- Human Control of Technology
- Safety and Security
- Transparency and Explainability
- Accountability
- Privacy

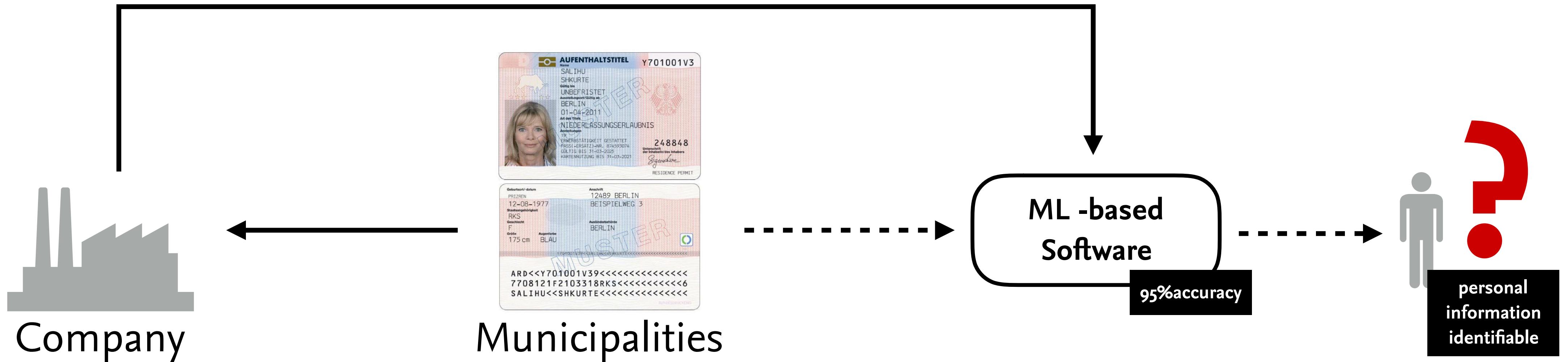




Accountability



Use Case: Automatic Anonymization



‘Algorithmic Accountability’ is the assessment of algorithms to discern and track bias, presuppositions, and prejudices built into, or resulting from algorithms affecting the society.

Wieringa, M. (2020). What to account for when accounting for algorithms (Vol. 40, pp. 1–18). Presented at the FAT* '20: Conference on Fairness, Accountability, and Transparency, New York, NY, USA: ACM.



Defining Accountability

“

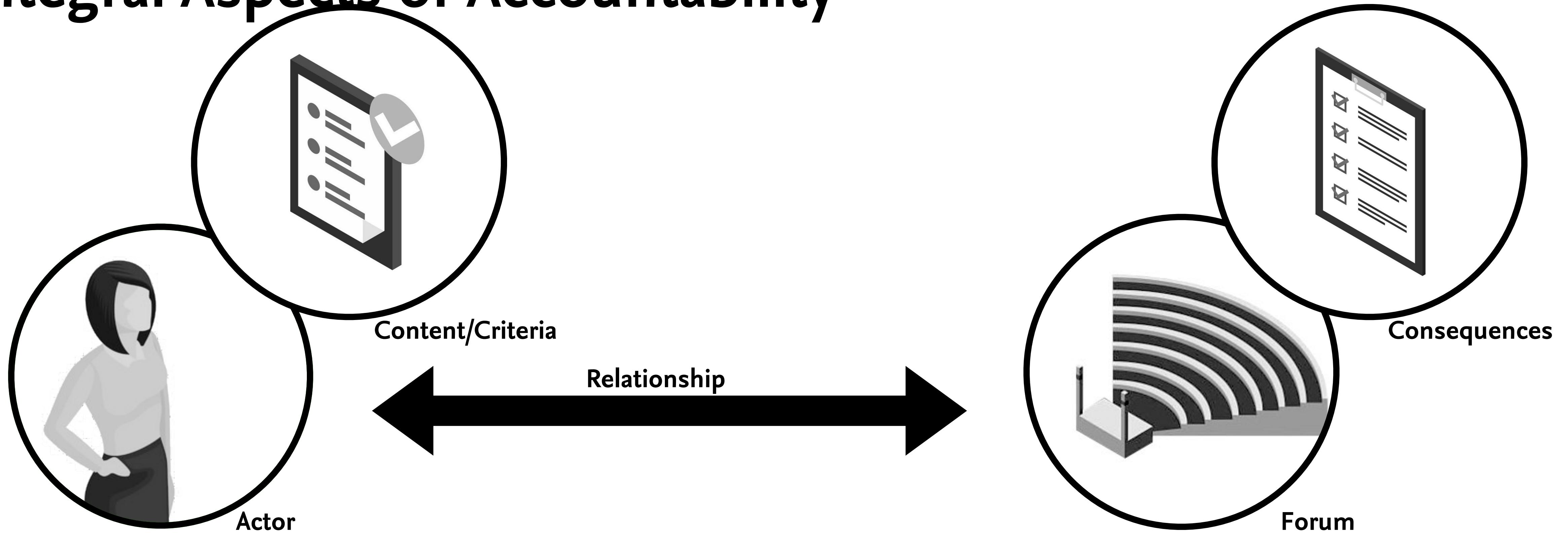
Accountability is: “a relationship between an actor and a forum, in which the actor has an obligation to explain and to justify his or her conduct, the forum can pose questions and pass judgement, and the actor may face consequences.”

Mark Bovens. (2007) Analysing and Assessing Accountability: A Conceptual Framework. European Law Journal 13, 4 (2007), 447–468.



Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn

Integral Aspects of Accountability



Mark Bovens. (2007) Analysing and Assessing Accountability: A Conceptual Framework. European Law Journal 13, 4 (2007), 447–468.



Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn

Accountability Relations by Forum



Mark Bovens. (2007) Analysing and Assessing Accountability: A Conceptual Framework. European Law Journal 13, 4 (2007), 447–468.



Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn



Photo by Kenny Eliason on Unsplash

Connectivity

The Growing Problem of Bots That Fight Online

The way software agents interact on the Web is poorly understood. Now evidence shows that they fight each other for years.

by Emerging Technology from the arXiv September 20, 2016



Technology

Wikipedia bots spent years fighting silent, tiny battles with each other

And no one even noticed

By Sara Chodosh February 27, 2017



Battle of the Bots: 'Main Reason for Conflicts is Lack of Central Supervision'

become a supporter subscribe / find a job

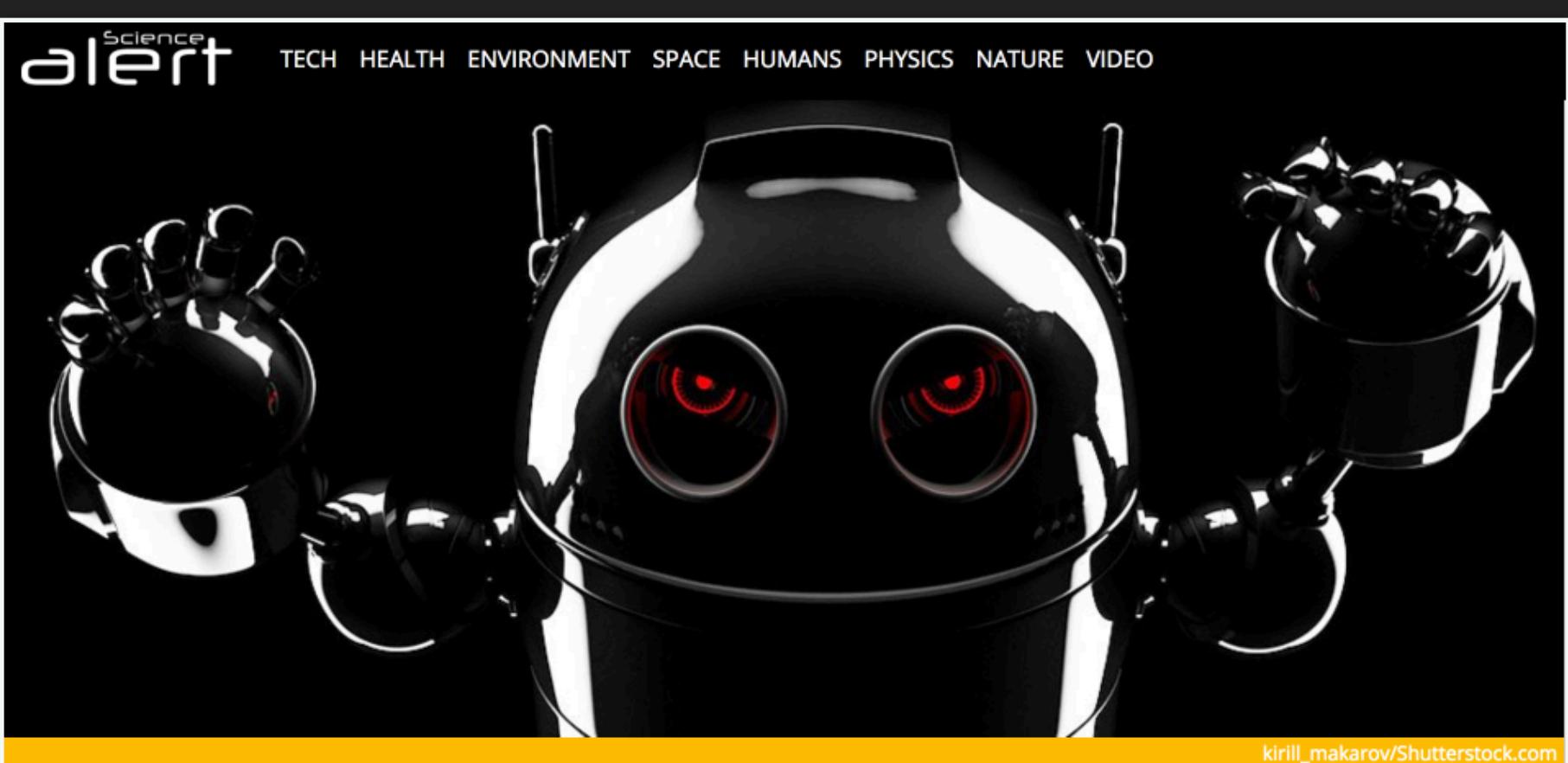
the guardian

news / opinion / sport / arts / life

tech / world / UK / science / cities / global development / more

Wikipedia

Study reveals bot-on-bot editing wars raging on Wikipedia's pages



kirill_makarov/Shutterstock.com

Investigation Reveals That Wikipedia's Bots Are in a Silent, Never-Ending War With Each Other

Technology can be brutal.

PETER DOCKRILL 26 FEB 2017

Home › 'Computer bots are like humans, having fights lasting years'

'Computer bots are like humans, having fights lasting years'

INNOVATION

Automated Wikipedia Edit-Bots Have Been Fighting Each Other For A Decade

In 13 different languages.

MATT SIMON BUSINESS 03.01.17 07:00 AM

INTERNET BOTS FIGHT EACH OTHER BECAUSE THEY'RE ALL TOO HUMAN



Use Case: Even Good Bots Fight



The screenshot shows the PLOS ONE article page for the paper "Even good bots fight: The case of Wikipedia". The page includes the PLOS ONE logo, navigation links (PUBLISH, ABOUT, BROWSE, SEARCH), and user account options (create account, sign in). The article title is "Even good bots fight: The case of Wikipedia" by Milena Tsvetkova, Ruth García-Gavilanes, Luciano Floridi, Taha Yasseri. It was published on February 23, 2017. The metrics section shows 105 saves, 34 citations, 54,874 views, and 314 shares. Below the metrics are download, print, and share buttons, and a "Check for updates" link. The abstract section discusses the increasing number of bots online and their interactions, mentioning that while they are intended to support Wikipedia, they often undo each other's edits. The subject areas listed include Online encyclopedias, Language, Internet, Network reciprocity, Twitter, Artificial intelligence, Encyclopedias, and Ecosystems.

Abstract

In recent years, there has been a huge increase in the number of bots online [...]. The online world has turned into an ecosystem of bots. However, our knowledge of how these automated agents are interacting with each other is rather poor. [...] In this article, we analyze the interactions between bots that edit articles on Wikipedia. We track the extent to which bots undid each other's edits over the period 2001–2010, model how pairs of bots interact over time, and identify different types of interaction trajectories. We find that, although Wikipedia bots are intended to support the encyclopedia, they often undo each other's edits and these sterile “fights” may sometimes continue for years.[...]

IS THERE A REPRODUCIBILITY CRISIS?



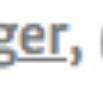
©nature

Replication Study by Geiger & Halfaker

RESEARCH-ARTICLE OPEN ACCESS

Operationalizing Conflict and Cooperation between Automated Software Agents in Wikipedia: A Replication and Expansion of 'Even Good Bots Fight'

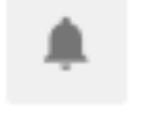
    

Authors:  [R. Stuart Geiger](#),  [Aaron Halfaker](#) [Authors Info & Affiliations](#)

Publication: Proceedings of the ACM on Human-Computer Interaction • December 2017 • Article No.: 49

- <https://doi.org/10.1145/3134684>

R. Stuart Geiger and Aaron Halfaker. 2017. Operationalizing Conflict and Cooperation between Automated Software Agents in Wikipedia: A Replication and Expansion of 'Even Good Bots Fight'. Proc. ACM Hum.-Comput. Interact. 1, CSCW, Article 49 (November 2017), 33 pages. DOI:<https://doi.org/10.1145/3134684>



What is Reproducibility?

“

A research project is computationally reproducible if a second investigator (including you in the future) can recreate the final reported results of the project, including key quantitative findings, tables, and figures, given only a set of files and written instructions.

Kitzes, J., Turek, D., & Deniz, F. (Eds.). (2018). *The Practice of Reproducible Research: Case Studies and Lessons from the Data-Intensive Sciences*. Oakland, CA: University of California Press. <http://www.practicereproducibleresearch.org/>



Repeatability vs. Reproducibility vs. Replicability

Repeatability — same team, same experimental setup

» [...] For computational experiments, this means that a researcher can reliably repeat her own computation.

Reproducibility — different team, same experimental setup

» [...] For computational experiments, this means that an independent group can obtain the same result using the author's own artifacts.

Replicability — different team, different experimental setup

» [...] For computational experiments, this means that an independent group can obtain the same result using artifacts which they develop completely independently.

Degree of Replication

Same Measurement and Analysis

Same Data Set

1. Checking the analysis

Different Measurement and Analysis

Same Population

3. Exact replication

2. Re-analyzing the data

4. Conceptual extension

Different Population

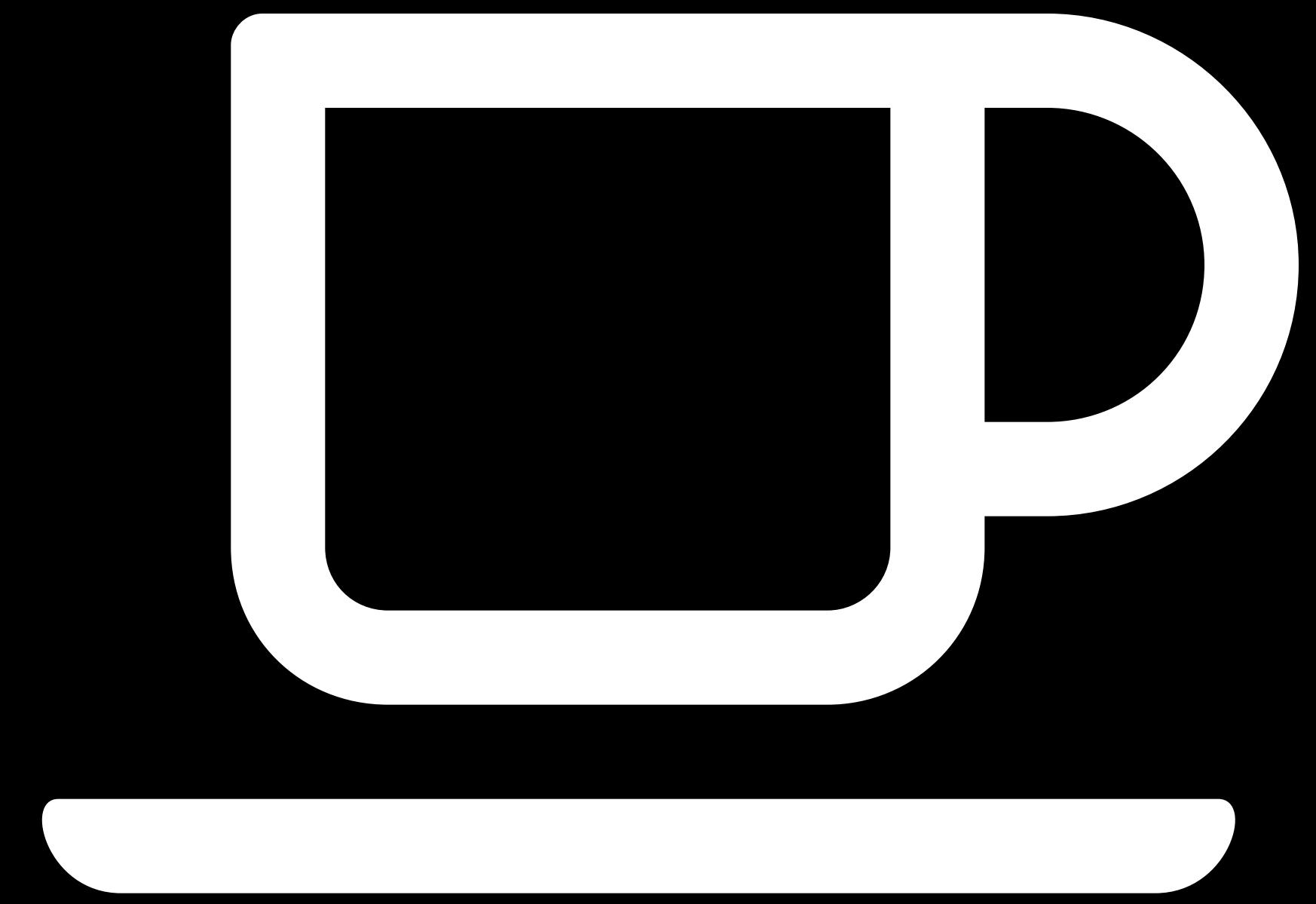
5. Empirical generalization

6. Generalization and extension

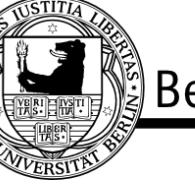


Translating this Principle in your Data Science Practice





5 minutes break



Translating this Principle in your Data Science Practice



Why is it useful to think about Reproducibility?

Publishing your research openly shows that you take responsibility for your research. Including the possibility that you might be wrong.

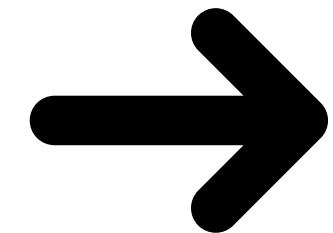
It provides transparency around your values, motivations, and assumptions.

Having a public audience in mind when designing and publishing your projects encourages you to reflect on your values, motivations, assumptions, and thought process—and how that might influence your project.

Taking a Human-Centered Approach

Audience Who are you publishing your research for?

Values



Purpose How do you want them to use your research?

Context What factors (under your control) will impact whether/how they use it?

Thinking in terms of “human-centeredness” can help you think of trade-offs of opening your data. For example, it can help you decide when/what NOT to publish openly!

What should you Consider in Your Data Science Practice?

1. **Support using your data**

Make your collected data freely available by considering the principles of FAIR data.

2. **Support sharing your research**

License your data/scripts/software.

3. **Support understanding your research**

Prepare your methods and the entire process of data analysis as reproducible as possible and provide relevant documentation.

Principles of FAIR Data

Findable

Metadata and data should be easy to find for both humans and computers. Machine-readable metadata are essential for automatic discovery of datasets and services.

Interoperable

The data usually need to be integrated with other data. In addition, the data need to interoperate with applications or workflows for analysis, storage, and processing.

Accessible

Once the user finds the required data, she/he needs to know how they can be accessed, possibly including authentication and authorization.

Reusable

To achieve this, metadata and data should be well-described so that they can be replicated and/or combined in different settings.

Upload Your Data with a DOI

Making your project identifiable via a Digital Object Identifiers (DOIs)

- Unique ID for an artefact
- Works even if sites fall over
- Single point of reference



Comment: Why uploading data, when people can ask you?

100% of the authors in the following studies signed to share the data upon request.

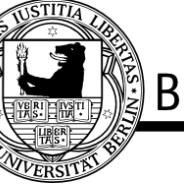
Actual sharing rate

- » **27%** out of 141 requests (Wicherts et al., 2006)
- » **38%** out of 394 requests (Vanpaemel et al., 2015)
- » **44%** out of 204 requests provided some „artifacts“ (Stodden et al., 2018), where 26% could be reproduced.

„Sharing upon request“ does not work.

Wicherts, J. M., Borsboom, D., Kats, J., & Molenaar, D. (2006). <http://doi.org/10.1037/0003-066X.61.7.726> | Vanpaemel, W., Vermorgen, M., Deriemaeker, L., & Storms, G. (2015). <http://doi.org/10.1525/collabra.13>
Stodden, V., Seiler, J., & Ma, Z. (2018). <http://doi.org/10.1073/pnas.1708290115> | Slide adapted from Felix Schönbrodt “Changing incentive structures to foster the actual sharing rate of open data” (2018)





What should you consider in your data science practice?

1. Support using your data

Make your collected data freely available by considering the principles of FAIR data.

2. Support sharing your research

License your data/scripts/software.

3. Support understanding your research

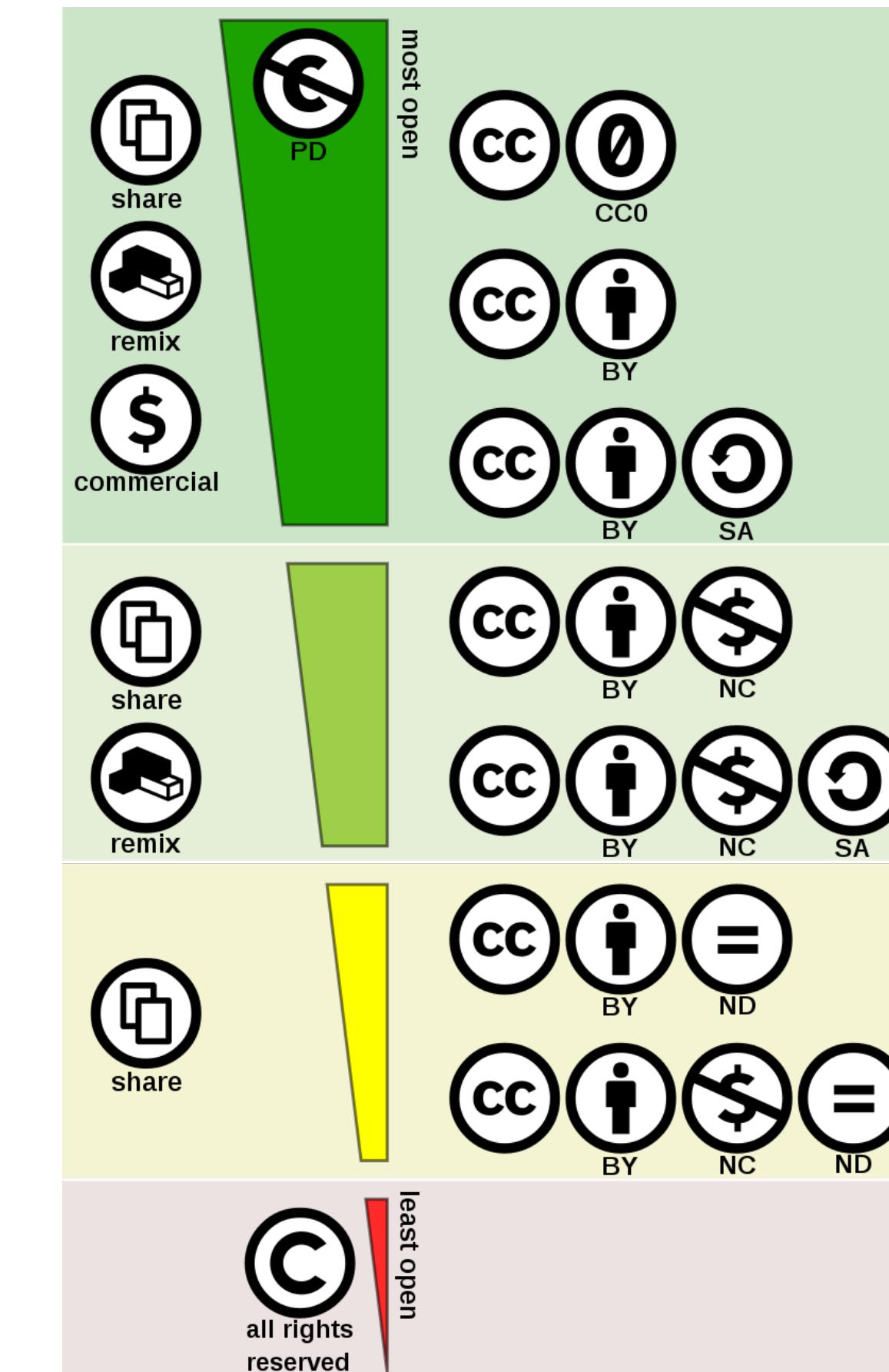
Prepare your methods and the entire process of data analysis as reproducible as possible and provide relevant documentation.



License Your Data and Documentation

For **licensing data or documentation**, one of the Creative Commons (CC) licenses can be employed.

It allows the creator of the data/documentation, to give other people the right to share, use, and build upon the data/documentation.



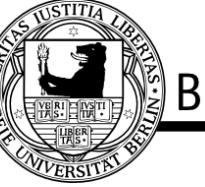
Common building blocks:

- » 0: public domain release
- » BY: must provide attribution
- » NC: cannot use commercially
- » ND: cannot make derivatives
- » SA: must release derivatives under the same/a compatible license

Image by Shaddim, CC BY 4.0, <https://commons.wikimedia.org/w/index.php?curid=47247325>

Description from https://en.wikipedia.org/w/index.php?title=Creative_Commons_license&oldid=984371920





License Your Software

A **software is a legal instrument** governing the use or redistribution of software.

There are **many different open source licenses** available such as the MIT license or the General Public License v3.0. But there are quite different, while the MIT License lets people do almost anything they want with your source code, the GPL v3.0 is much stricter.

As a rule of thumb, If you extend existing code, use the license of the community in which you are participating.

Choose an open source license

An open source license protects contributors and users. Businesses and savvy developers won't touch a project without this protection.

{ Which of the following best describes your situation? }



I need to work in a community.



I want it simple and permissive.



I care about sharing improvements.

Use the [license preferred by the community](#) you're contributing to or the [MIT License](#) is short and to the point. It lets people do almost anything they want with your project, like making and distributing closed source versions. The [GNU GPLv3](#) also lets people do almost anything they want with your project, except distributing closed source versions. Description from https://en.wikipedia.org/w/index.php?title=Software_license&oldid=983929529 <https://choosealicense.com/licenses/>





Example

Instructional Material

All Software Carpentry and Data Carpentry instructional material is made available under the [Creative Commons Attribution license](#). The following is a human-readable summary of (and not a substitute for) the [full legal text of the CC BY 4.0 license](#).

You are free:

- to **Share**—copy and redistribute the material in any medium or format
- to **Adapt**—remix, transform, and build upon the material

for any purpose, even commercially.

The licensor cannot revoke these freedoms as long as you follow the license terms.

Under the following terms:

- **Attribution**—You must give appropriate credit (mentioning that your work is derived from work that is Copyright © Software Carpentry and, where practical, linking to <http://software-carpentry.org/>), provide a [link to the license](#), and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.

No additional restrictions—You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits. With the understanding that:

Notices:

- You do not have to comply with the license for elements of the material in the public domain or where your use is permitted by an applicable exception or limitation.
- No warranties are given. The license may not give you all of the permissions necessary for your intended use. For example, other rights such as publicity, privacy, or moral rights may limit how you use the material.

Software

Except where otherwise noted, the example programs and other software provided by Software Carpentry and Data Carpentry are made available under the [OSI-approved MIT license](#).

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

Trademark

"Software Carpentry" and "Data Carpentry" and their respective logos are registered trademarks of [NumFOCUS](#).

<https://reproducible-science-curriculum.github.io/sharing-RR-Jupyter/LICENSE.html>

Copyright © 2016–2019 Data Carpentry

[Edit on GitHub](#) / [Contributing](#) / [Source](#) / [Cite](#) / [Contact](#)





What should you consider in your data science practice?

1. Support using your data

Make your collected data freely available by considering the principles of FAIR data.

2. Support sharing your research

License your data/scripts/software.

3. Support understanding your research

Prepare your methods and the entire process of data analysis as reproducible as possible and provide relevant documentation.

Slide adapted from Jonathan T. Morgan & Oliver Keyes (Course DATA 512)



Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn

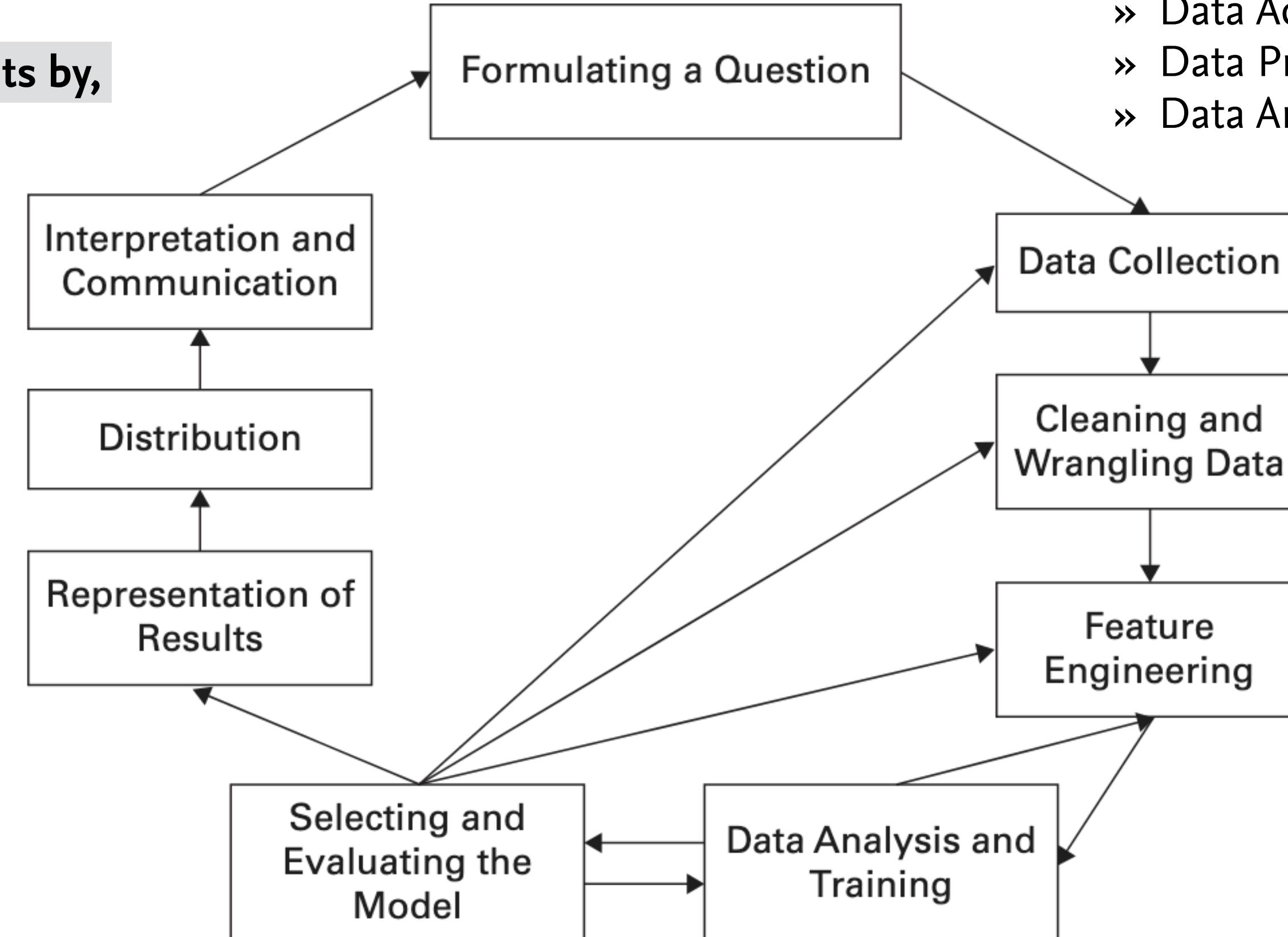
Stages of a Basic Reproducible Workflow

**Analyze and derive improvements by,
e.g.**

- » Bug tracker
- » Interviews
- » Checklists

**Open workflow for feedback
of, e.g.**

- » Peers
- » Colleagues
- » Community



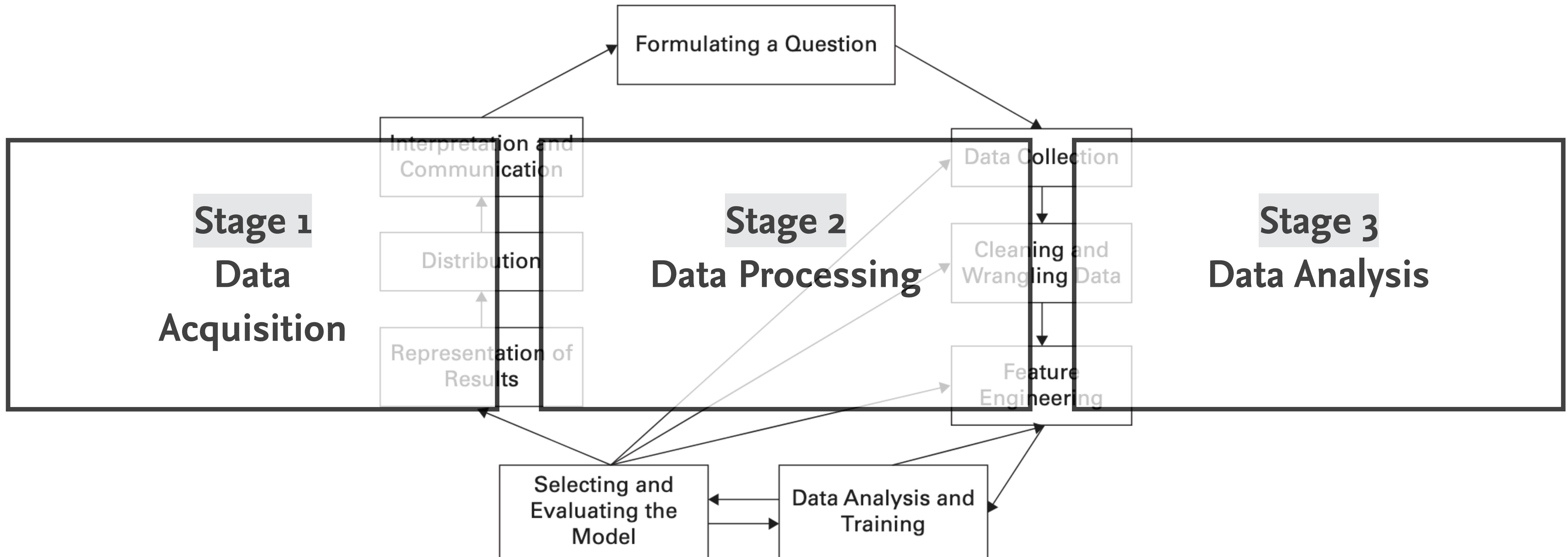
Re-consider the stages of a reproducible workflow

- » Data Acquisition
- » Data Processing
- » Data Analysis

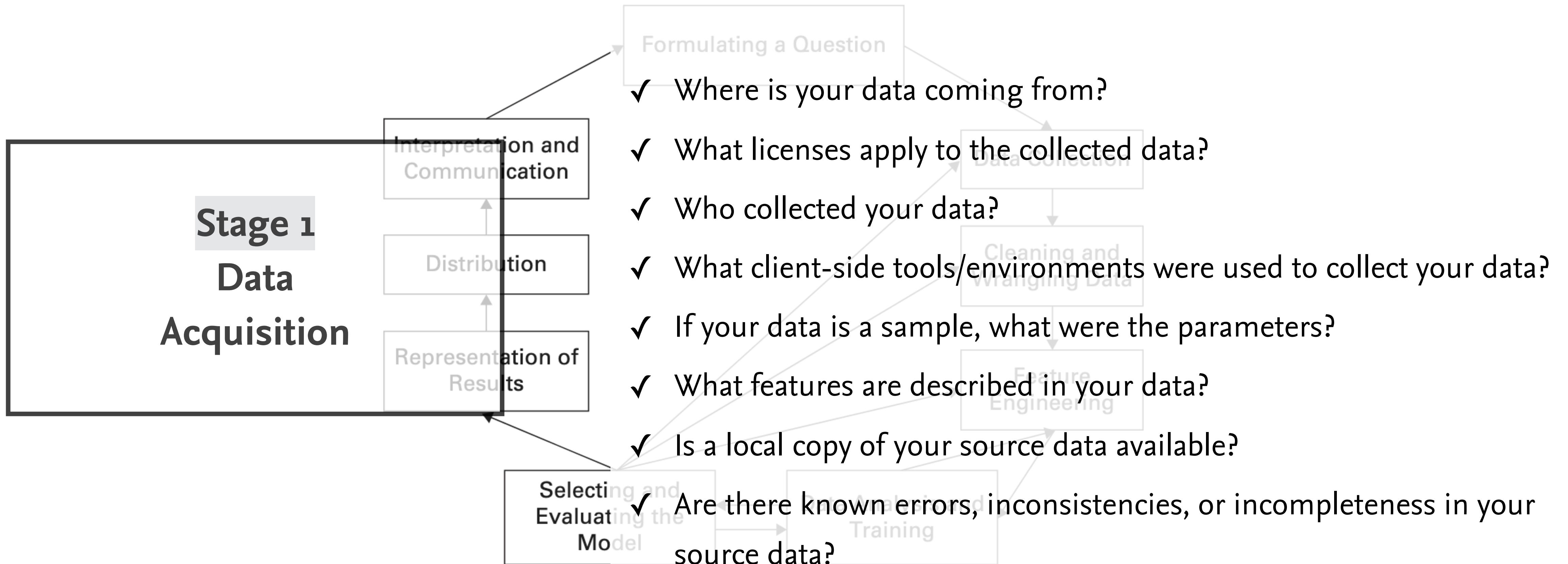
**Build new/Refactor existing workflow
by, e.g.**

- » Defining data origin and license
- » Documenting used libraries
- » Providing test data

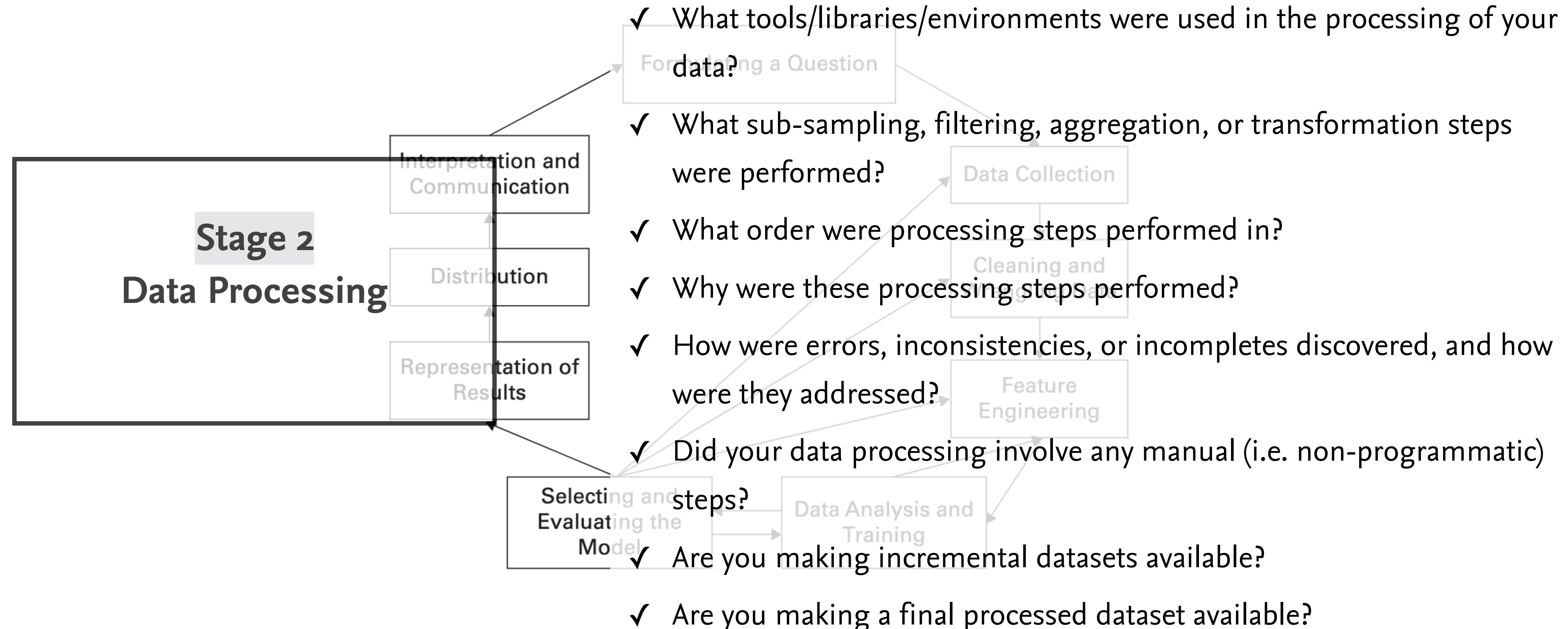
Stages of the Basic Reproducible Workflow



Stages of the Basic Reproducible Workflow



Stages of the Basic Reproducible Workflow

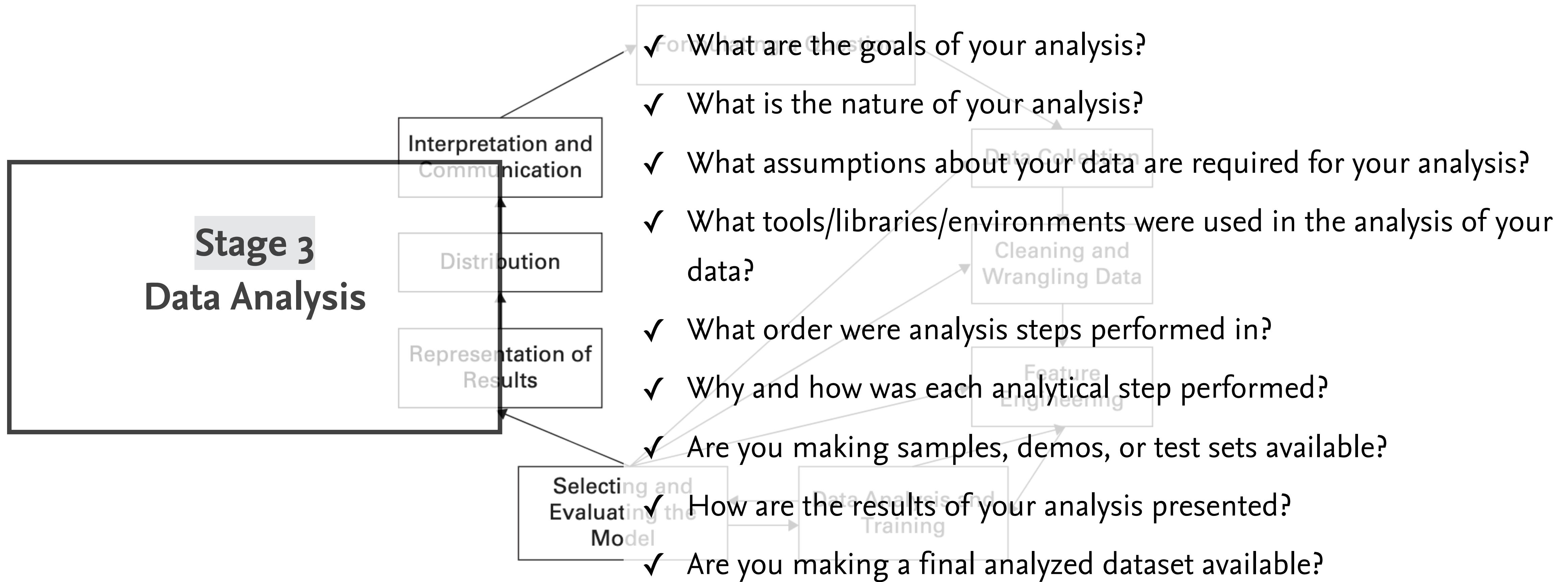


Slide adapted from Jonathan T. Morgan & Oliver Keyes (Course DATA 512)



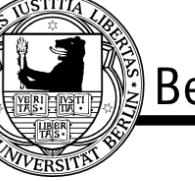
Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn

Stages of the Basic Reproducible Workflow



Three Key Practices for your Reproducible Workflow

1. Clearly separate, label, and document all data, files, and operations that occur on data and files.
2. Document all operations fully, automating them as much as possible, and avoiding manual intervention in the workflow when feasible.
3. Design a workflow as a sequence of small steps that are glued together, with intermediate outputs from one step feeding into the next step as inputs.



Use Case (continued)



Article and Data provided by Tsvetkova et al.

Copyright: © 2017 Tsvetkova et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability: All relevant data are available from figshare at [10.6084/m9.figshare.4597918](https://doi.org/10.6084/m9.figshare.4597918).



Even Good Bots Fight: The Case of Wikipedia

Cite Download (3.47 GB) Share Embed + Collect

Dataset posted on 31.01.2017 by Milena Tsvetkova

Each file contains the history of all the articles in the WP. Articles are separated by their names within the file. Each line of the file below the name of an article, contains a delimiter "^^^" followed by the timestamp of each edit, a binary flag of 0/1 corresponding to a normal/revert edit, an accenting integer code, starting from 1, assigned to each new revision, whose text is not similar to any of the previous ones, otherwise the same code as the previous version with the similar text, and finally the editor of the version.

FUNDING

HORIZON 2020 NO 645043

HISTORY

31.01.2017 - First online date, Posted date

figshare

788 views | 462 downloads | 1 citations

CATEGORIES

- Social and Community Informatics

KEYWORD(S)

Wikipedia | editorial wars | social network

LICENCE

CC BY 4.0

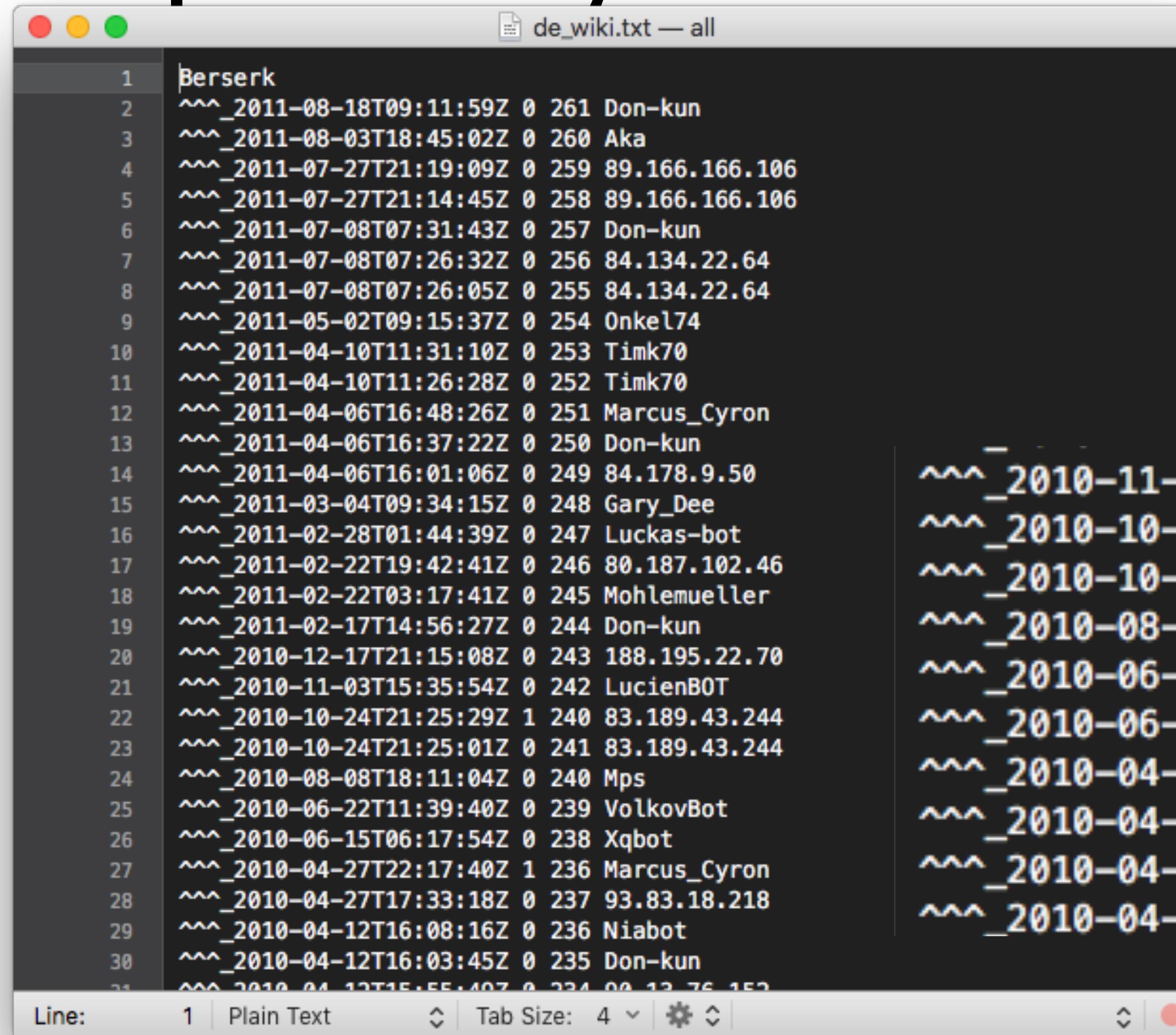
EXPORTS

Tsvetkova, M., García-Gavilanes, R., Floridi, L., & Yasseri, T. (2017). Even good bots fight: The case of Wikipedia. PLoS ONE, 12(2).





Data provided by Tsvetkova et al. article



```
de_wiki.txt — all
1 Berserk
2 ^^^_2011-08-18T09:11:59Z 0 261 Don-kun
3 ^^^_2011-08-03T18:45:02Z 0 260 Aka
4 ^^^_2011-07-27T21:19:09Z 0 259 89.166.166.106
5 ^^^_2011-07-27T21:14:45Z 0 258 89.166.166.106
6 ^^^_2011-07-08T07:31:43Z 0 257 Don-kun
7 ^^^_2011-07-08T07:26:32Z 0 256 84.134.22.64
8 ^^^_2011-07-08T07:26:05Z 0 255 84.134.22.64
9 ^^^_2011-05-02T09:15:37Z 0 254 Onkel74
10 ^^^_2011-04-10T11:31:10Z 0 253 Timk70
11 ^^^_2011-04-10T11:26:28Z 0 252 Timk70
12 ^^^_2011-04-06T16:48:26Z 0 251 Marcus_Cyron
13 ^^^_2011-04-06T16:37:22Z 0 250 Don-kun
14 ^^^_2011-04-06T16:01:06Z 0 249 84.178.9.50
15 ^^^_2011-03-04T09:34:15Z 0 248 Gary_Dee
16 ^^^_2011-02-28T01:44:39Z 0 247 Luckas-bot
17 ^^^_2011-02-22T19:42:41Z 0 246 80.187.102.46
18 ^^^_2011-02-22T03:17:41Z 0 245 Mohlemueller
19 ^^^_2011-02-17T14:56:27Z 0 244 Don-kun
20 ^^^_2010-12-17T21:15:08Z 0 243 188.195.22.70
21 ^^^_2010-11-03T15:35:54Z 0 242 LucienBOT
22 ^^^_2010-10-24T21:25:29Z 1 240 83.189.43.244
23 ^^^_2010-10-24T21:25:01Z 0 241 83.189.43.244
24 ^^^_2010-08-08T18:11:04Z 0 240 Mps
25 ^^^_2010-06-22T11:39:40Z 0 239 VolkovBot
26 ^^^_2010-06-15T06:17:54Z 0 238 Xqbot
27 ^^^_2010-04-27T22:17:40Z 1 236 Marcus_Cyron
28 ^^^_2010-04-27T17:33:18Z 0 237 93.83.18.218
29 ^^^_2010-04-12T16:08:16Z 0 236 Niabot
30 ^^^_2010-04-12T16:03:45Z 0 235 Don-kun
31 ^^^_2010-04-12T15:55:40Z 0 234 00.12.76.152
```

Line: 1 | Plain Text | Tab Size: 4 | 

The data processing step is already incompletely documented. Missing information are the user groups (user, bots).
Are vandals excluded?

```
^__ 2010-11-03T15:35:54Z 0 242 LucienBOT
^__ 2010-10-24T21:25:29Z 1 240 83.189.43.244
^__ 2010-10-24T21:25:01Z 0 241 83.189.43.244
^__ 2010-08-08T18:11:04Z 0 240 Mps
^__ 2010-06-22T11:39:40Z 0 239 VolkovBot
^__ 2010-06-15T06:17:54Z 0 238 Xqbot
^__ 2010-04-27T22:17:40Z 1 236 Marcus_Cyron
^__ 2010-04-27T17:33:18Z 0 237 93.83.18.218
^__ 2010-04-12T16:08:16Z 0 236 Niabot
^__ 2010-04-12T16:03:45Z 0 235 Don-kun
```



Replication Study by Geiger & Halfaker

RESEARCH-ARTICLE OPEN ACCESS

Operationalizing Conflict and Cooperation between Automated Software Agents in Wikipedia: A Replication and Expansion of 'Even Good Bots Fight'



Authors:  R. Stuart Geiger,  Aaron Halfaker [Authors Info & Affiliations](#)

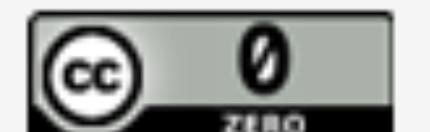
Publication: Proceedings of the ACM on Human-Computer Interaction • December 2017 • Article No.: 49

• <https://doi.org/10.1145/3134684>

5 424



LICENCE



CC0



Code, data, and Jupyter notebooks available at: <https://github.com/halfak/are-the-bots-really-fighting>
<https://doi.org/10.6084/m9.figshare.5362216>



R. Stuart Geiger and Aaron Halfaker. 2017. Operationalizing Conflict and Cooperation between Automated Software Agents in Wikipedia: A Replication and Expansion of 'Even Good Bots Fight'. Proc. ACM Hum.-Comput. Interact. 1, CSCW, Article 49 (November 2017), 33 pages. DOI:<https://doi.org/10.1145/3134684>



Course «Human-Centered Data Science» | Summer Term 2022 | Claudia Müller-Birn

Data provided by Geiger & Halfaker

halfak / [are-the-bots-really-fighting](#)

Code Issues Pull requests Actions Projects Security Insights

Join GitHub today

GitHub is home to over 50 million developers working together to host and review code, manage projects, and build software together.

Dismiss

Sign up

master ▾ 13 branches 3 tags Go to file Code ▾

staeiou add doi	263bd96 on Nov 7, 2018	228 commits
analysis	add docs and fix path for data files	3 years ago
datasets	update datasets after rerun	3 years ago
docs	move tables from anon repo	3 years ago
environment	tweak description	3 years ago
paper	Add files via upload	3 years ago
sql	Fixes issues in bot_revert_monthly and adds monthly stats to ma...	4 years ago

About

A research project exploring revert patterns between bots on Wikipedia.

Readme

Releases 3

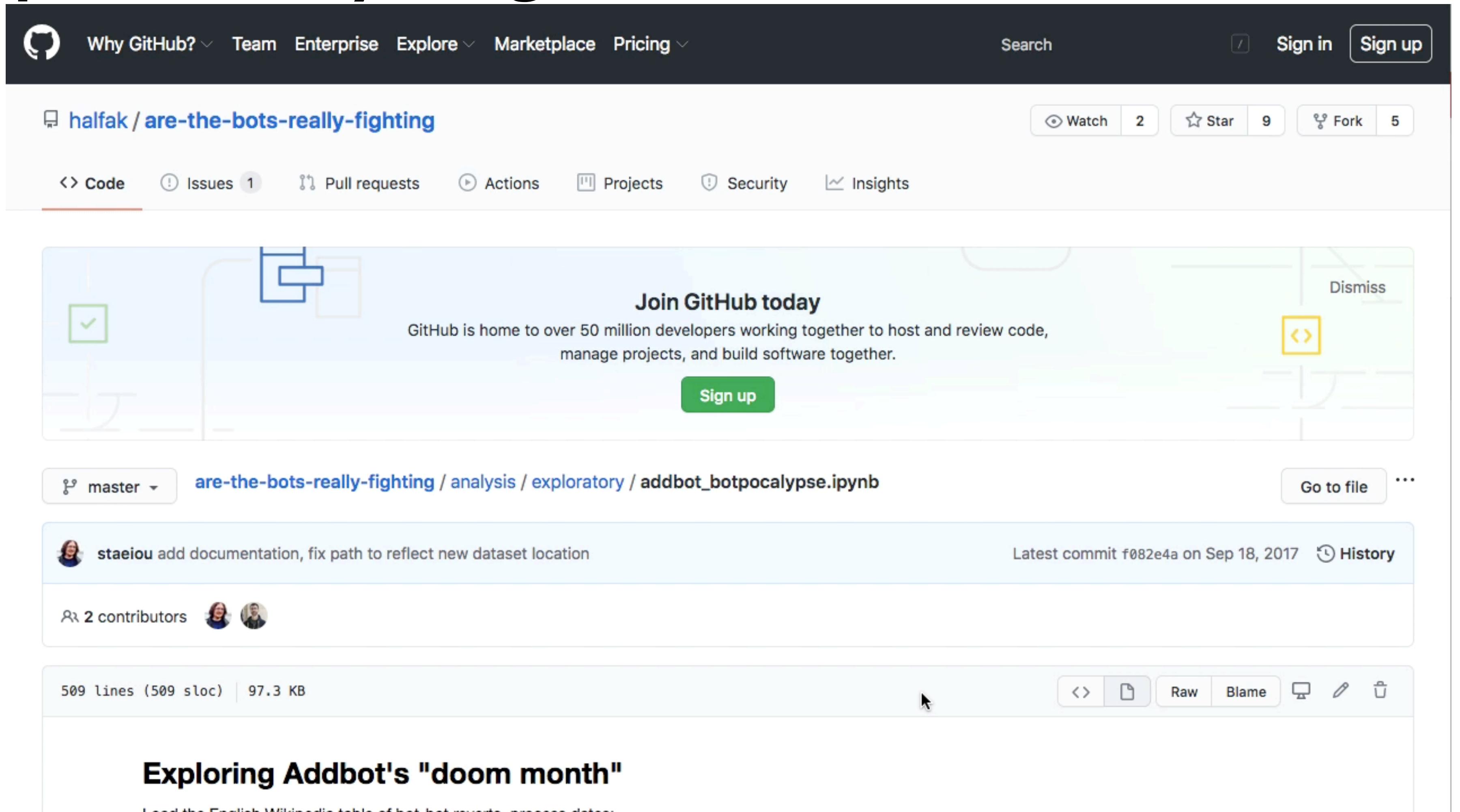
Photo ready Latest on Sep 9, 2017

+ 2 releases

<https://github.com/halfak/are-the-bots-really-fighting/>



Scripts provided by Geiger & Halfaker



The screenshot shows a GitHub repository page for the user 'halfak' with the repository name 'are-the-bots-really-fighting'. The page includes a navigation bar with links to 'Why GitHub?', 'Team', 'Enterprise', 'Explore', 'Marketplace', and 'Pricing'. A search bar and sign-in links are also present. Below the header, there are buttons for 'Watch' (2), 'Star' (9), and 'Fork' (5). The main content area displays a large image of a Jupyter notebook titled 'addbot_botpocalypse.ipynb'. A modal window titled 'Join GitHub today' is overlaid on the page, encouraging users to sign up. The notebook file details show it was last updated on Sep 18, 2017, by 'staeiou' with a commit message: 'add documentation, fix path to reflect new dataset location'. It has 2 contributors. The file size is 509 lines (509 sloc) and 97.3 KB. Below the file details, the title 'Exploring Addbot's "doom month"' is visible, along with a note: 'Load the English Wikipedia table of hot-bot reverts... process data'.



What are the results of Geiger & Halfaker's replication?

They could show that the overwhelming majority of bot-bot reverts constitute routine, productive, and even collaborative work between bots.

They define bots as assemblages of “code and a human developer” which are responsible for operating the bot in alignment with Wikipedia’s complex policy environment.

**Why these researchers arrived at this very different perspective,
we will discuss in another lecture.**

Check your Insights

- » How do Accountability and Reproducibility relate?
- » What is the difference between repeatability, reproducibility, replicability?
- » How can you differentiate various forms of replication studies?
- » What aspects should you consider in your (open) data science practice?
- » What does it mean - FAIR data?
- » What is Creative Commons (CC)?
- » What should be your key practice of a basic reproducible workflow?





«Human-Centered Data Science»

Next week: Examining the Role of Data in Your Data Science Practice

Prof. Dr. Claudia Müller-Birn

Human-Centered Computing, Institute of Computer Science

Freie Universität Berlin

May 19, 2022