

things go
better
with
Cake

Small Mistake, Huge Difference.
Ensure Data Quality.



7 4



Data Quality in the Age of AI

CIRES PhD School 2024

Based on research with Hazar Harmouch, Lisa Ehrlinger, Sedir Mohammed, and Divesh Srivastava

October 29, 2024

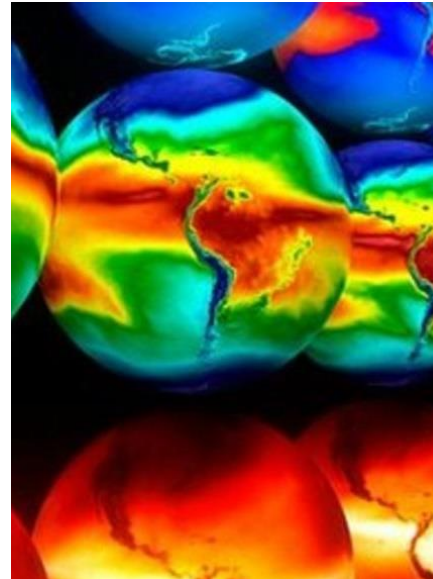
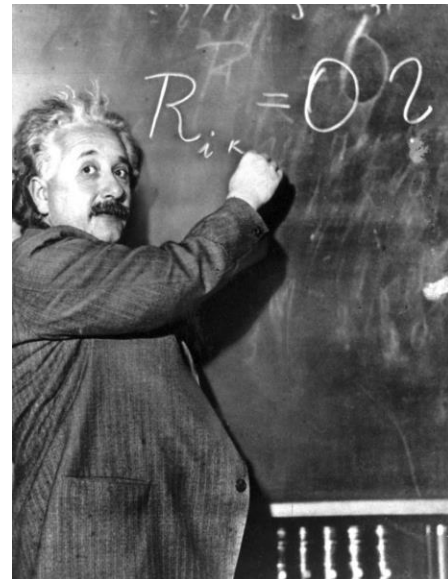
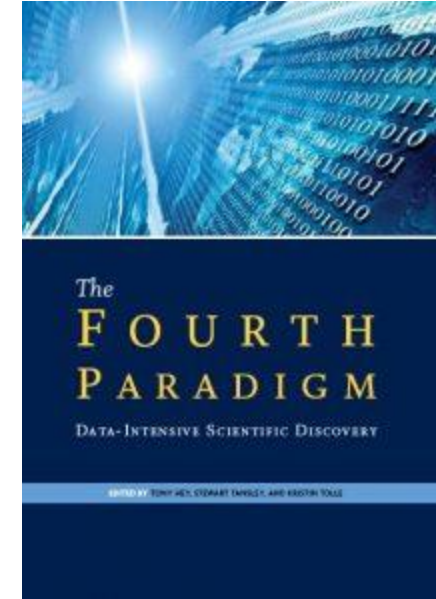
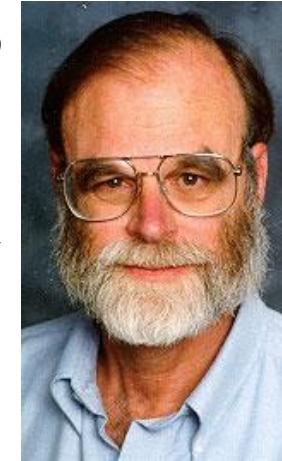
Felix Naumann

The Fourth Paradigm of Science

1. Empirical and experimental
2. Theoretical
3. Computational
4. Data-intensive
5. Intelligence-driven and knowledge-centric

We have to do better producing tools to *support the whole research cycle* – from data capture and data curation to data analysis and data visualization.

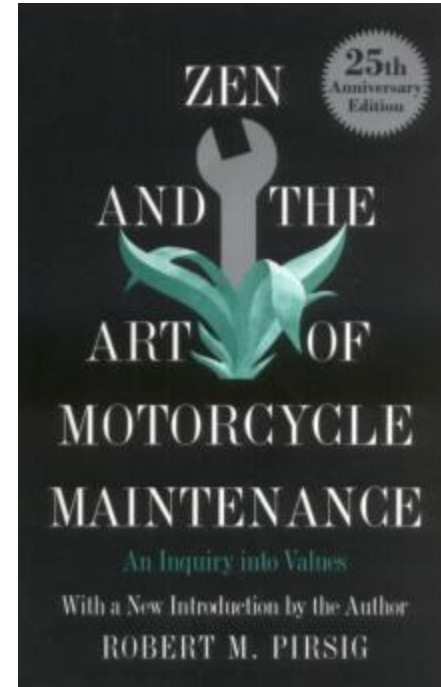
Jim Gray



Felix Naumann
Data Quality

*“Even though quality
cannot be defined, you
know what it is.”*

Robert Pirsig




Felix Naumann
Data Quality

Data Errors by Data Quality Researchers




Jana BAUCKMANN
*Hasso-Plattner-Institut, University of
Potsdam*
GERMANY



Alexander ALBRECHT
*Hasso-Plattner-Institut, Universität
Potsdam*
GERMANY



Christoph BÖHM
Hasso-Plattner-Institut, Potsdam
GERMANY





Frank KAUFER
*Hasso Plattner Institute, Potsdam
University*
GERMANY



Felix NAUMANN
Hasso Plattner Institute
GERMANY

Felix Naumann
Data Quality

Vandalism in Wikipedia Tables

| No. | | Mayor | Took Office | Left Office | Prior Experience | Deputy Mayor |
|-----|--|--------------|------------------|-------------------|--|---------------|
| 62 | | Mel Lastman | January 1, 1998 | November 30, 2003 | Mayor of North York (1969–1997) | Case Ootes |
| 63 |  | David Miller | December 1, 2003 | November 30, 2010 | City Councillor for Parkdale-High Park (1994–2003) | Joe Pantalone |
| 64 |  | Rob | | | | |

| | |
|--|--|
| - I —Non-Hispanic II 31.7% II 37.9% II 59.0% <ref name="fifteen">From 15% sample</ref> II 91.2% | + I —Non-Hispanic II 20.7% II 21.9% II 59.0% <ref name="fifteen">From 15% sample</ref> II 91.2% |
| - I [[African American Black or African American]] II 32.9% II 39.1% II 32.7% II 8.2% | + I [[African American Black or African American]] II 50.9% II 49.1% II 42.7% II 8.2% |

Example for vandalism in Wikipedia tables: Tampering with the proportions of ethnic minorities.
 [https://en.wikipedia.org/w/index.php?title=Chicago&diff=prev&oldid=654893961]

Hidden Values / Hidden Value

| Feld | Name1 | Name2 | Name3 | City | District | Street | Sum |
|-------------------|--------|-------|-------|------|----------|--------|--------|
| Mobile phone | 41 | 501 | 10 | 0 | 2677 | 297 | 3526 |
| Phone | 15 | 98 | 6 | 0 | 221 | 9579 | 9919 |
| Cost center | 283 | 1112 | 73 | 2 | 87 | 16 | 1573 |
| Registration ID | 11 | 583 | 1 | 1 | 0 | 3 | 599 |
| Delivery ID | 55 | 390 | 9 | 0 | 212 | 15 | 681 |
| Department | 3711 | 9997 | 115 | 60 | 439 | 175 | 14497 |
| Embargo flag | 129 | 143 | 2 | 0 | 66 | 9 | 349 |
| Deletion flag | 1028 | 442 | 5 | 36 | 113 | 10 | 1634 |
| Legal form | 131700 | 66136 | 187 | 6 | 64 | 57 | 198150 |
| Credit info | 0 | 100 | 11 | 0 | 18 | 0 | 129 |
| Commission | 216 | 352 | 1 | 2 | 36 | 10 | 617 |
| Construction site | 2013 | 3452 | 42 | 5 | 124 | 222 | 5858 |
| Loading point | 2923 | 3808 | 94 | 1503 | 958 | 3065 | 12351 |
| Administration | 13410 | 12461 | 172 | 19 | 295 | 7075 | 33432 |
| Summe | 155535 | 99575 | 728 | 1634 | 5310 | 20533 | |

Felix Naumann
Data Quality

DQ-Problems: Effects

- Incorrect prices in inventory retail databases
 - Costs for consumers 2.5 billion \$
 - 80% of barcode-scan-errors to the disadvantage of consumer
- IRS 1992: almost 100,000 tax refunds not deliverable
- 50% to 80% of computerized criminal records in the U.S. were found to be inaccurate, incomplete, or ambiguous.
- US-Postal Service: of 100,000 mass-mailings up to 7,000 undeliverable due to incorrect addresses
- Poor AI system performance

**IRS might
be after you
— to mail
you a check**

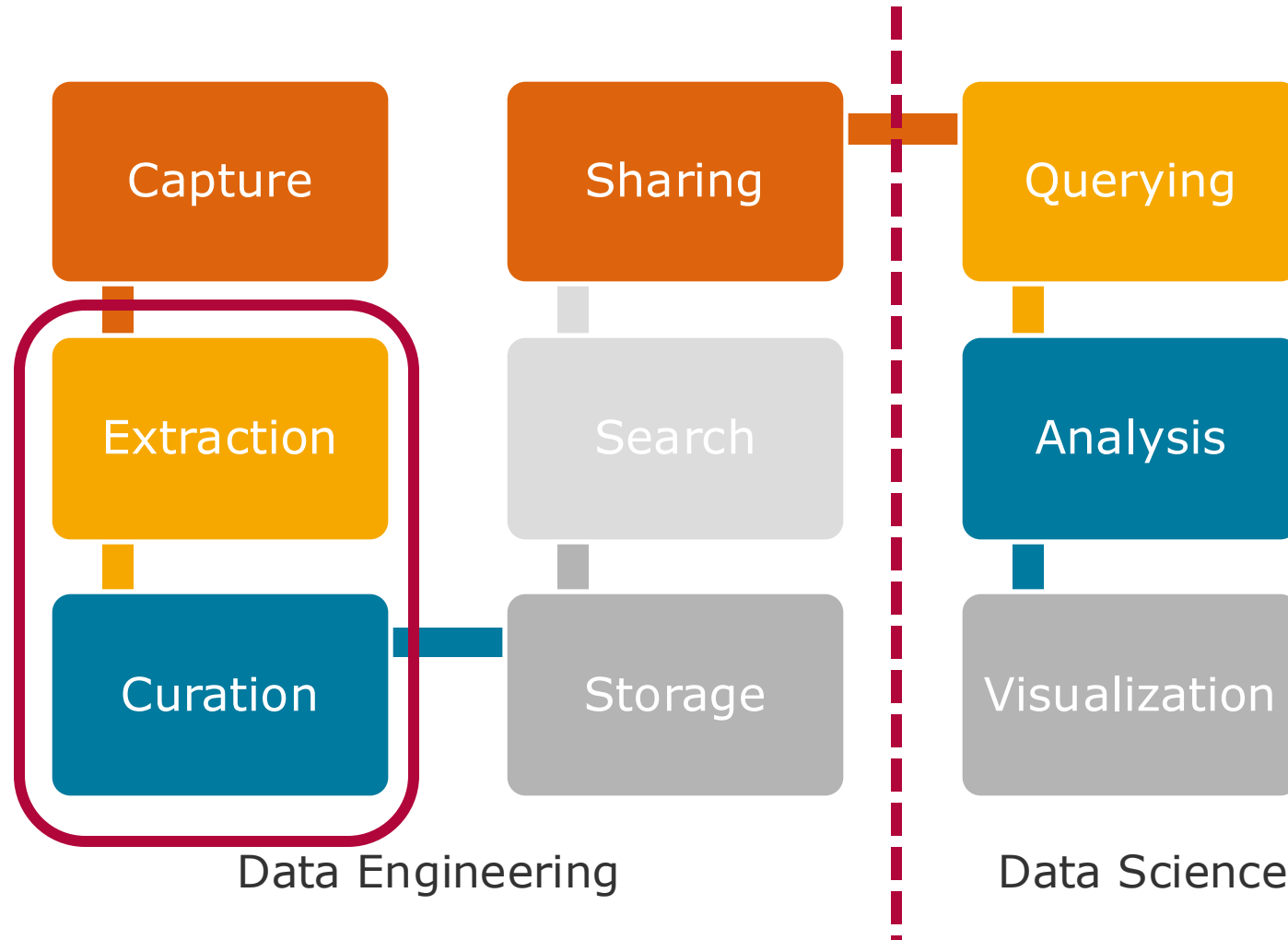
Incorrect addresses
stall nearly 1,500
Tennessee refunds

By **BONNA de la CRUZ**
Staff Writer

Now that Tilcia L. Menifee knows that she'll be getting \$500 in a tax refund from Uncle Sam, she can do some Christmas shopping, she said.

Felix Naumann
Data Quality

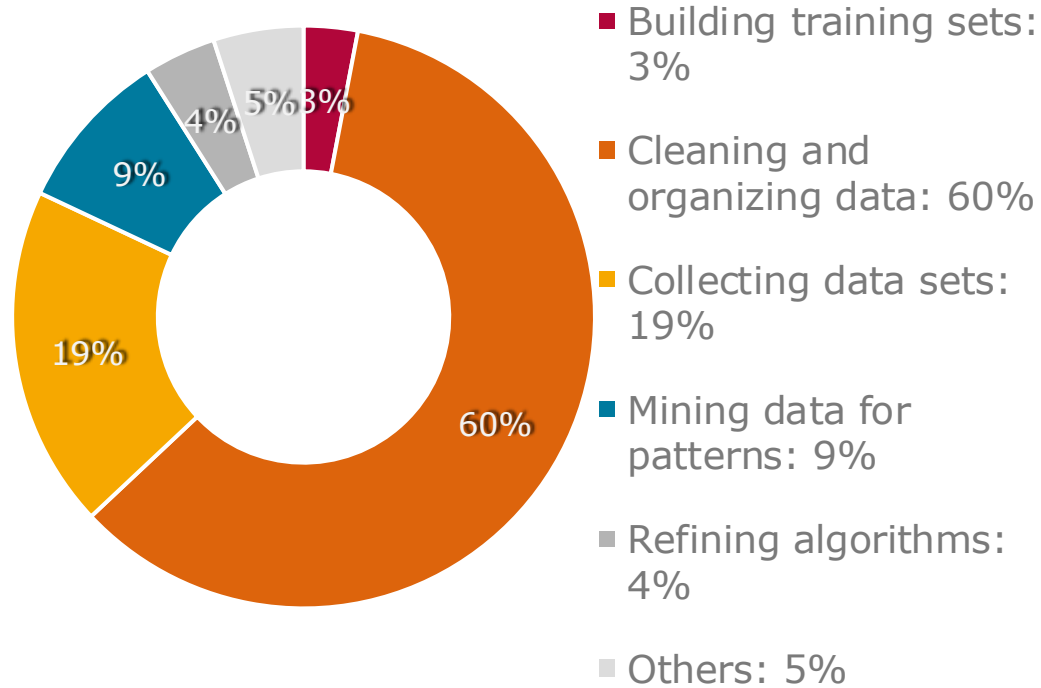
Data Science Pipeline



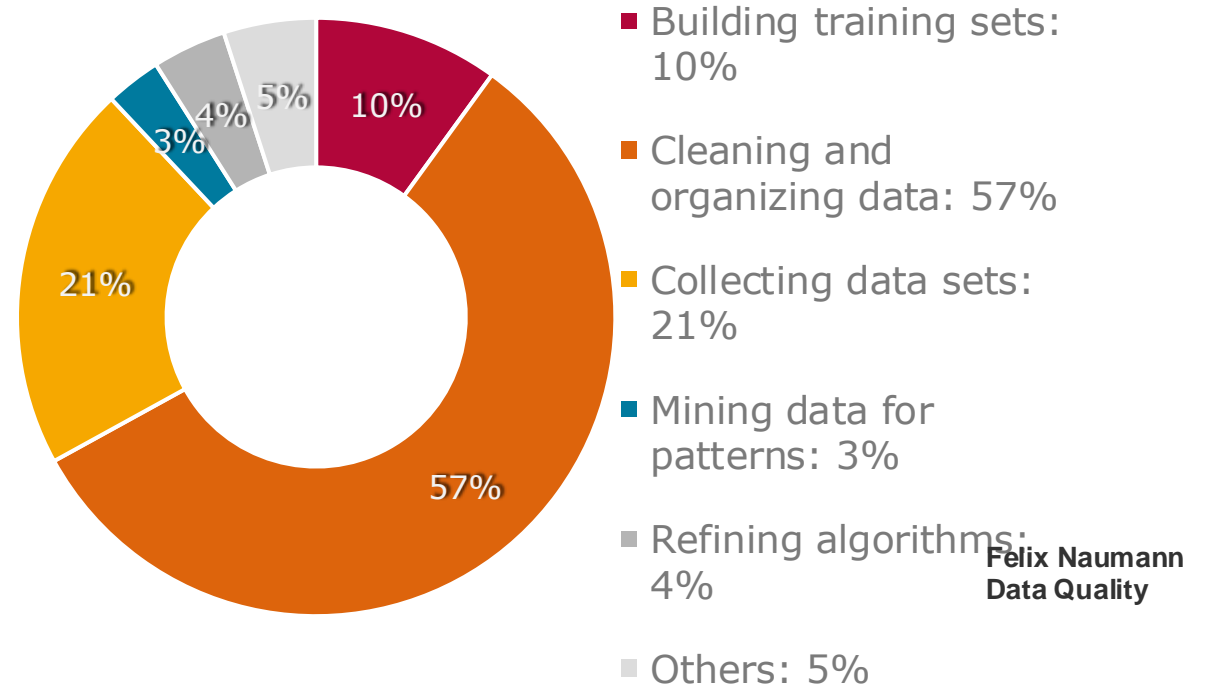
Felix Naumann
Data Quality

Data preparation in reality

What data scientists spend the **most time** doing?



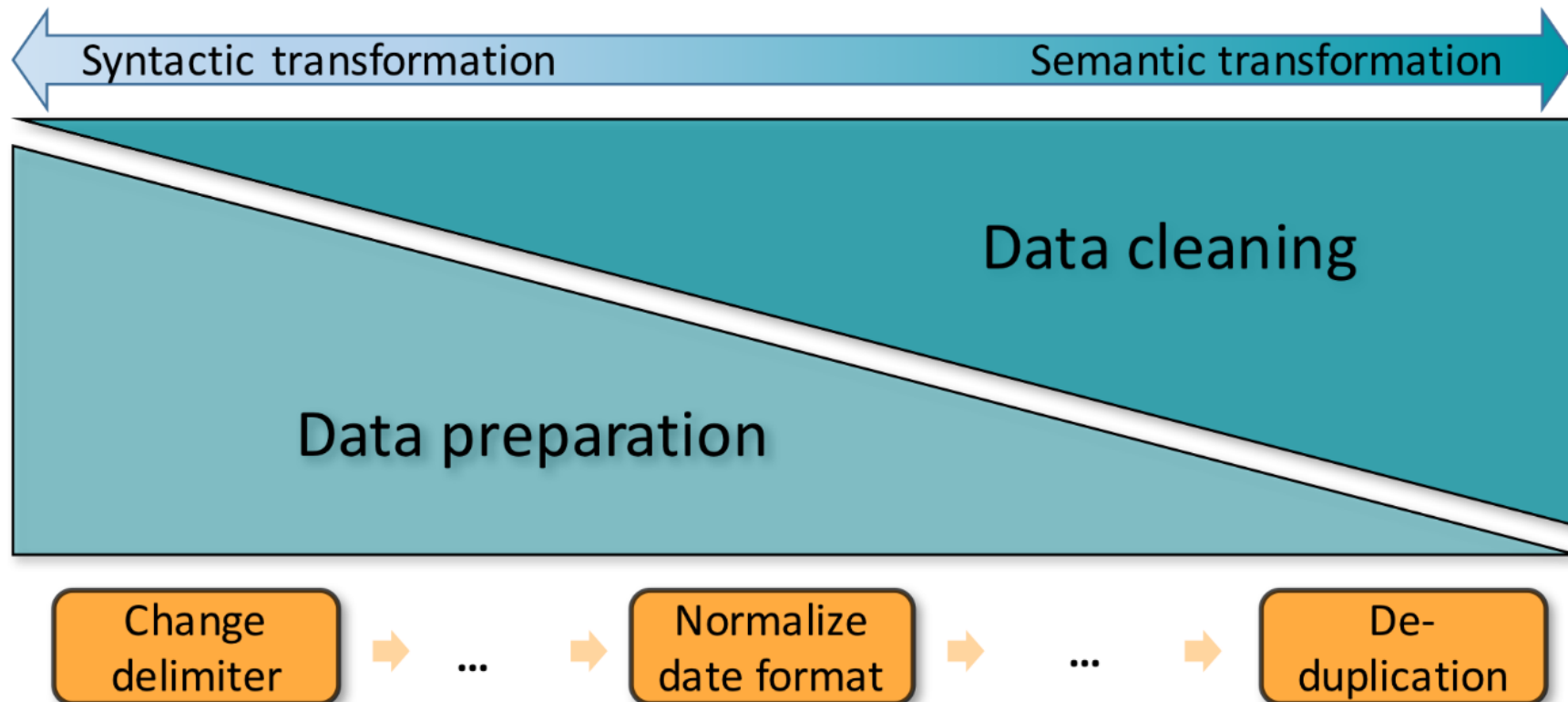
What is the **least enjoyable** part of data science?



Felix Naumann
Data Quality

Data Preparation vs. Data Cleaning

- Data preparation adds syntactic and structural value
- Data cleaning adds semantic value



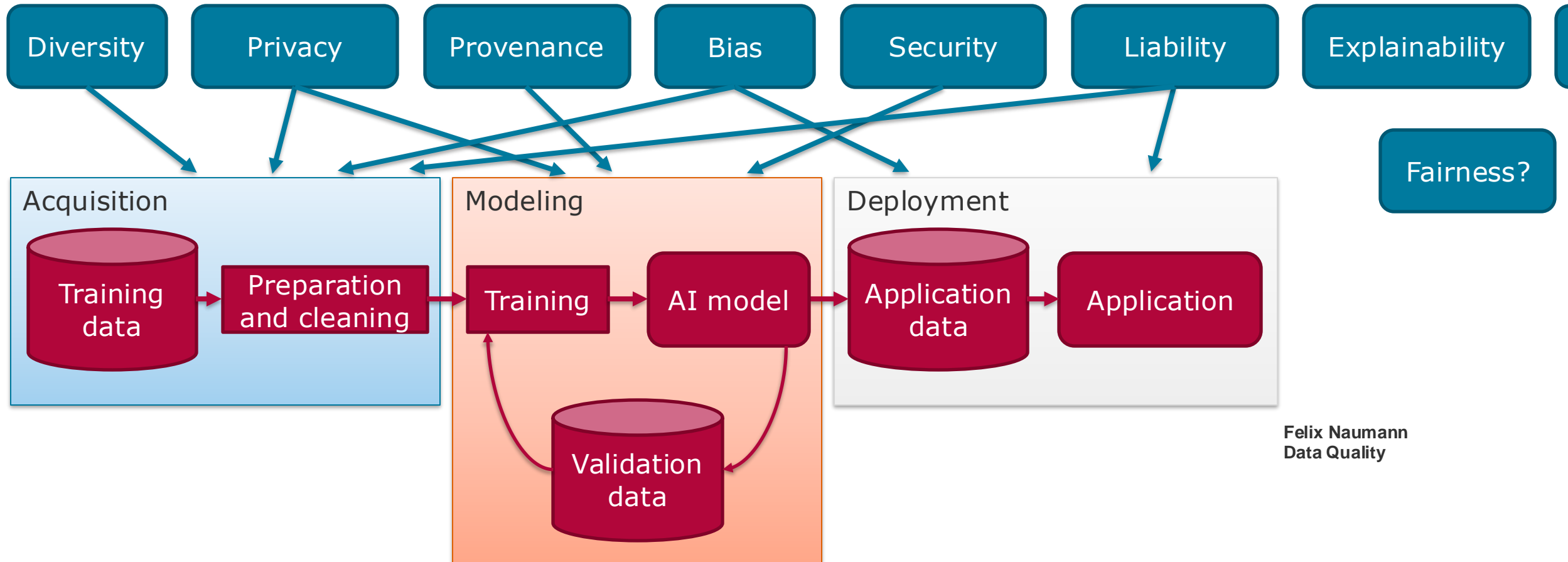
Felix Naumann
Data Quality

Agenda

- 1. Data and Information Quality Research**
2. Data Preparation
3. Data Quality and AI Systems
4. Data Quality Assessment

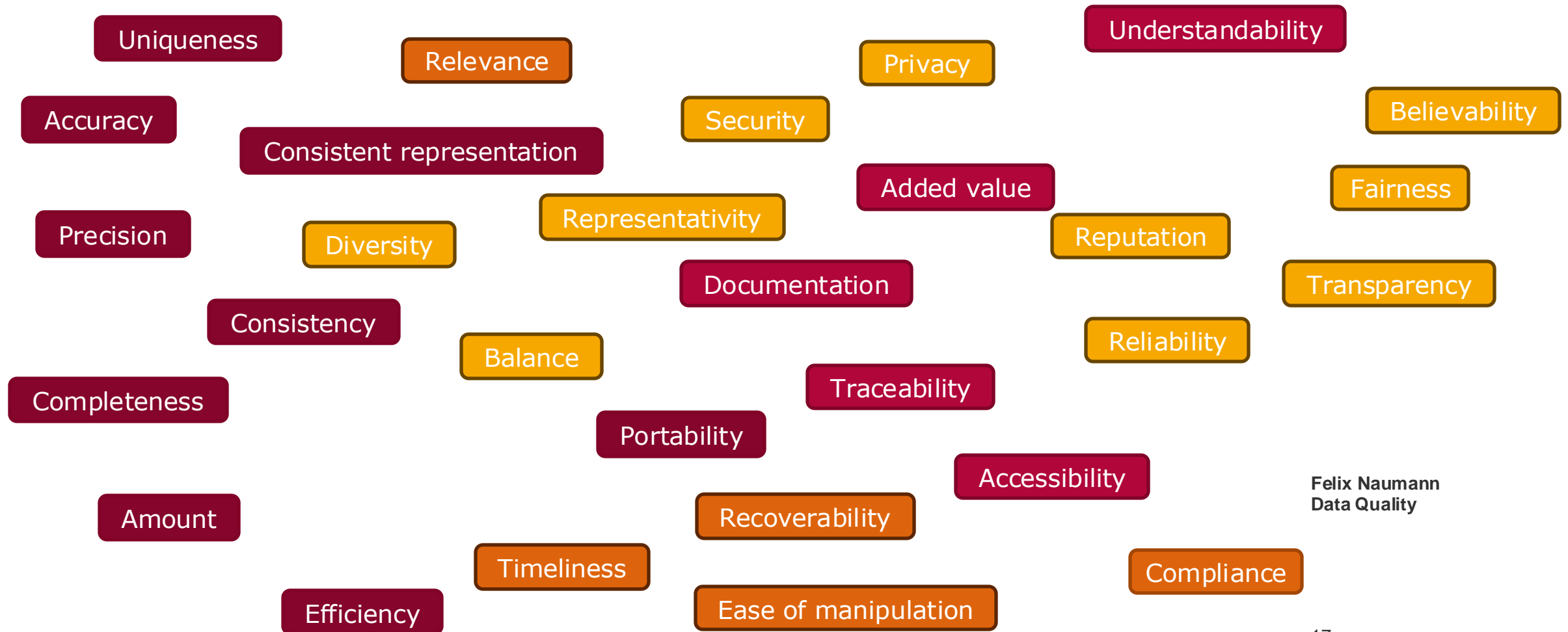


New AI-specific Data Quality Dimensions



Felix Naumann
Data Quality

28 DQ Dimensions



Felix Naumann
Data Quality

Agenda

1. Data and Information Quality Research
- 2. Data Preparation**
3. Data Quality and AI Systems
4. Data Quality Assessment



| ID | Name | Sex | Age | Height | Weight | Team | NOC | Games | Year | Season | City | Sport | Event | Medal |
|----|--|-----|-----|----------|------------------|-------|---------------|-------|----------|------------------|------------------------|---|----------|-------|
| 1 | "A Dijkstra" | "M" | 24 | 180,80 | "China" | "CHN" | "1992 Summer" | 1992 | "Summer" | "Barcelona" | "Basketball" | "Basketball Men's Basketball" | NA | |
| 2 | "A Lamusi" | "M" | 23 | 170,60 | "China" | "CHN" | "2012 Summer" | 2012 | "Summer" | "London" | "Judo" | "Judo Men's Extra-Lightweight" | NA | |
| 3 | "Gunnar Nielsen Aaby" | "M" | 24 | NA,NA | "Denmark" | "DEN" | "1920 Summer" | 1920 | "Summer" | "Antwerpen" | "Football" | "Football Men's Football" | NA | |
| 4 | "Edgar Lindenau Aaby" | "M" | 34 | NA,NA | "Denmark/Sweden" | "DEN" | "1900 Summer" | 1900 | "Summer" | "Paris" | "Tug-Of-War" | "Tug-Of-War Men's Tug-Of-War" | "Gold" | |
| 5 | "Christine Jacoba Aaftink" | "F" | 21 | 185,82 | "Netherlands" | "NED" | "1988 Winter" | 1988 | "Winter" | "Calgary" | "Speed Skating" | "Speed Skating Women's 500 metres" | NA | |
| 6 | "Christine Jacoba Aaftink" | "F" | 21 | 185,82 | "Netherlands" | "NED" | "1988 Winter" | 1988 | "Winter" | "Calgary" | "Speed Skating" | "Speed Skating Women's 1,000 metres" | NA | |
| 7 | "Christine Jacoba Aaftink" | "F" | 25 | 185,82 | "Netherlands" | "NED" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Speed Skating" | "Speed Skating Women's 500 metres" | NA | |
| 8 | "Christine Jacoba Aaftink" | "F" | 25 | 185,82 | "Netherlands" | "NED" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Speed Skating" | "Speed Skating Women's 1,000 metres" | NA | |
| 9 | "Christine Jacoba Aaftink" | "F" | 27 | 185,82 | "Netherlands" | "NED" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Speed Skating" | "Speed Skating Women's 500 metres" | NA | |
| 10 | "Christine Jacoba Aaftink" | "F" | 27 | 185,82 | "Netherlands" | "NED" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Speed Skating" | "Speed Skating Women's 1,000 metres" | NA | |
| 11 | "Per Knut Aaland" | "M" | 31 | 188,75 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 10 kilometres" | NA | |
| 12 | "Per Knut Aaland" | "M" | 31 | 188,75 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 50 kilometres" | NA | |
| 13 | "Per Knut Aaland" | "M" | 31 | 188,75 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 10/15 kilometres Pursuit" | NA | |
| 14 | "Per Knut Aaland" | "M" | 31 | 188,75 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 4 x 10 kilometres Relay" | NA | |
| 15 | "Per Knut Aaland" | "M" | 33 | 188,75 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 10 kilometres" | NA | |
| 16 | "Per Knut Aaland" | "M" | 33 | 188,75 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 30 kilometres" | NA | |
| 17 | "Per Knut Aaland" | "M" | 33 | 188,75 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 10/15 kilometres Pursuit" | NA | |
| 18 | "Per Knut Aaland" | "M" | 33 | 188,75 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 4 x 10 kilometres Relay" | NA | |
| 19 | "John Aalberg" | "M" | 31 | 183,72 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 10 kilometres" | NA | |
| 20 | "John Aalberg" | "M" | 31 | 183,72 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 50 kilometres" | NA | |
| 21 | "John Aalberg" | "M" | 31 | 183,72 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 10/15 kilometres Pursuit" | NA | |
| 22 | "John Aalberg" | "M" | 31 | 183,72 | "United States" | "USA" | "1992 Winter" | 1992 | "Winter" | "Albertville" | "Cross Country Skiing" | "Cross Country Skiing Men's 4 x 10 kilometres Relay" | NA | |
| 23 | "John Aalberg" | "M" | 33 | 183,72 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 10 kilometres" | NA | |
| 24 | "John Aalberg" | "M" | 33 | 183,72 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 30 kilometres" | NA | |
| 25 | "John Aalberg" | "M" | 33 | 183,72 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 10/15 kilometres Pursuit" | NA | |
| 26 | "John Aalberg" | "M" | 33 | 183,72 | "United States" | "USA" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Cross Country Skiing" | "Cross Country Skiing Men's 4 x 10 kilometres Relay" | NA | |
| 27 | "Cornelia ""Cor"" Aalten (-Strannood)" | "F" | 18 | 168 | "Netherlands" | "NED" | "1932 Summer" | 1932 | "Summer" | "Los Angeles" | "Athletics" | "Athletics Women's 100 metres" | NA | |
| 28 | "Cornelia ""Cor"" Aalten (-Strannood)" | "F" | 18 | 168 | "Netherlands" | "NED" | "1932 Summer" | 1932 | "Summer" | "Los Angeles" | "Athletics" | "Athletics Women's 4 x 100 metres Relay" | NA | |
| 29 | "Antti Sami Aalto" | "M" | 26 | 186,96 | "Finland" | "FIN" | "2002 Winter" | 2002 | "Winter" | "Salt Lake City" | "Ice Hockey" | "Ice Hockey Men's Ice Hockey" | NA | |
| 30 | "Einar Ferdinand ""Einari"" Aalto" | "M" | 26 | NA,NA | "Finland" | "FIN" | "1952 Summer" | 1952 | "Summer" | "Helsinki" | "Swimming" | "Swimming Men's 400 metres Freestyle" | NA | |
| 31 | "Jorma Ilmari Aalto" | "M" | 22 | 182,76.5 | "Finland" | "FIN" | "1980 Winter" | 1980 | "Winter" | "Lake Placid" | "Cross Country Skiing" | "Cross Country Skiing Men's 30 kilometres" | NA | |
| 32 | "Jyri Tapani Aalto" | "M" | 31 | 172,70 | "Finland" | "FIN" | "2000 Summer" | 2000 | "Summer" | "Sydney" | "Badminton" | "Badminton Men's Singles" | NA | |
| 33 | "Minna Maarit Aalto" | "F" | 30 | 159,55.5 | "Finland" | "FIN" | "1996 Summer" | 1996 | "Summer" | "Atlanta" | "Sailing" | "Sailing Women's Windsurfer" | NA | |
| 34 | "Minna Maarit Aalto" | "F" | 34 | 159,55.5 | "Finland" | "FIN" | "2000 Summer" | 2000 | "Summer" | "Sydney" | "Sailing" | "Sailing Women's Windsurfer" | NA | |
| 35 | "Pirjo Hannele Aalto (Mattila-)" | "F" | 32 | 171,65 | "Finland" | "FIN" | "1994 Winter" | 1994 | "Winter" | "Lillehammer" | "Biathlon" | "Biathlon Women's 7.5 kilometres Sprint" | NA | |
| 36 | "Arvo Ossian Aaltonen" | "M" | 22 | NA,NA | "Finland" | "FIN" | "1912 Summer" | 1912 | "Summer" | "Stockholm" | "Swimming" | "Swimming Men's 200 metres Breaststroke" | NA | |
| 37 | "Arvo Ossian Aaltonen" | "M" | 22 | NA,NA | "Finland" | "FIN" | "1912 Summer" | 1912 | "Summer" | "Stockholm" | "Swimming" | "Swimming Men's 400 metres Breaststroke" | NA | |
| 38 | "Arvo Ossian Aaltonen" | "M" | 30 | NA,NA | "Finland" | "FIN" | "1920 Summer" | 1920 | "Summer" | "Antwerpen" | "Swimming" | "Swimming Men's 200 metres Breaststroke" | "Bronze" | |
| 39 | "Arvo Ossian Aaltonen" | "M" | 30 | NA,NA | "Finland" | "FIN" | "1920 Summer" | 1920 | "Summer" | "Antwerpen" | "Swimming" | "Swimming Men's 400 metres Breaststroke" | "Bronze" | |
| 40 | "Arvo Ossian Aaltonen" | "M" | 34 | NA,NA | "Finland" | "FIN" | "1924 Summer" | 1924 | "Summer" | "Paris" | "Swimming" | "Swimming Men's 200 metres Breaststroke" | NA | |
| 41 | "Juhamatti Tapio Aaltonen" | "M" | 28 | 184,85 | "Finland" | "FIN" | "2014 Winter" | 2014 | "Winter" | "Sochi" | "Ice Hockey" | "Ice Hockey Men's Ice Hockey" | "Bronze" | |
| 42 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Individual All-Around" | "Bronze" | |
| 43 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Team All-Around" | "Gold" | |
| 44 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Floor Exercise" | NA | |
| 45 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Horse Vault" | "Gold" | |
| 46 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Parallel Bars" | NA | |
| 47 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Horizontal Bar" | NA | |
| 48 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Rings" | NA | |
| 49 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1948 Summer" | 1948 | "Summer" | "London" | "Gymnastics" | "Gymnastics Men's Pommel Horse" | "Gold" | |
| 50 | "Paavo Johannes Aaltonen" | "M" | 28 | 175,64 | "Finland" | "FIN" | "1952 Summer" | 1952 | "Summer" | "Helsinki" | "Gymnastics" | "Gymnastics Men's Individual All-Around" | NA | |

Data Preparation for AI: The Challenge

```

120 Nov-09,,4,47,35,17,99,32,1055,165,578,16,0,18,16,2,36,5,149,2,47,0,0,16,11,5,32,10,43,5,115,1
121 Dec-09,,3,41,32,15,89,27,930,145,566,14,0,17,17,2,36,4,131,2,49,0,0,12,10,5,27,8,40,6,106,1
122 Jan-10,,3,51,41,17,109,33,799,143,654,19,0,20,18,2,39,5,125,2,52,0,0,14,13,6,33,8,35,5,136,1
123 Feb-10,,3,46,36,14,96,32,636,133,545,17,0,19,15,1,35,4,97,1,44,0,0,13,12,6,31,8,24,4,113,1
124 Mar-10,,4,48,36,15,99,29,700,126,550,17,0,19,15,2,36,4,100,2,44,0,0,13,11,6,30,6,19,4,113,1
125 Apr-10,,4,57,42,19,119,33,792,157,665,20,0,24,17,3,44,4,115,2,52,0,0,17,15,8,39,7,21,5,141,1
126 May-10,,3,46,34,18,99,27,629,127,535,16,0,19,13,3,36,4,45,1,42,0,0,12,10,6,28,6,27,5,118,1
127 Jun-10,,3,43,33,20,97,26,682,132,531,14,0,18,13,5,36,4,55,1,39,0,0,11,10,6,29,6,27,5,115,1
128 Jul-10,,5,55,40,26,121,36,1075,182,662,Data are confidential,0,21,16,6,43,5,14,2,51,0,0,11,10,10,31,8,35,5,144,1
129 Aug-10,,5,43,32,20,95,28,987,165,553,Data are confidential,0,17,11,5,34,4,135,2,46,0,0,10,8,6,24,7,24,5,121,1
130 Sep-10,,7,48,34,18,100,33,957,158,562,Data are confidential,0,19,13,4,36,5,148,2,46,0,0,16,10,5,31,7,27,5,123,1
131 Oct-10,,9,63,44,22,129,49,1191,195,728,Data are confidential,0,24,19,4,47,6,197,3,57,0,0,1,22,13,6,41,10,29,7,157,1
132 Nov-10,,7,52,40,18,109,47,1047,183,605,Data are confidential,0,19,16,3,38,6,154,2,47,0,0,14,11,5,29,10,20,4,132,1
133 Dec-10,,6,55,42,18,114,41,1065,189,691,Data are confidential,0,21,20,3,43,5,167,3,54,0,0,0,14,11,6,31,8,20,4,143,1
134 Jan-11,,6,60,48,18,126,52,856,190,690,Data are confidential,0,22,20,3,45,6,149,2,52,0,0,1,16,15,7,38,10,19,4,157,1
135 Feb-11,,7,47,39,15,101,37,699,156,592,Data are confidential,0,19,16,2,37,4,115,2,46,0,0,14,12,5,32,8,13,2,123,1
136 Mar-11,,8,51,38,16,105,34,678,137,587,Data are confidential,0,20,16,2,37,4,115,2,46,0,0,14,12,5,32,8,13,2,123,1
137 Apr-11,,7,62,46,19,127,37,827,167,683,Data are confidential,0,23,18,4,45,5,118,2,60,0,0,15,12,5,32,7,15,3,143,0
138 May-11,,5,49,37,19,106,35,655,132,545,Data are confidential,0,19,14,4,36,5,49,2,45,0,0,11,10,6,27,7,17,3,122,0
139 Jun-11,,5,46,36,21,103,36,749,137,567,Data are confidential,0,17,13,5,35,5,72,2,45,0,0,10,8,6,25,8,21,2,127,0
140 Jul-11,,6,56,42,25,123,42,1133,189,728,Data are confidential,0,20,16,6,42,6,137,3,55,0,0,10,8,5,23,9,28,4,151,0
141 Aug-11,,5,45,34,18,97,34,956,153,594,Data are confidential,0,18,12,4,34,5,133,3,43,0,0,14,8,4,26,7,25,4,121,0
142 Sep-11,,7,51,36,17,104,40,992,153,621,Data are confidential,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,9,30,4,127,0
143 Oct-11,,8,61,45,18,125,53,1336,216,768,Data are confidential,0,22,20,2,45,8,191,3,68,0,1,20,11,5,36,12,34,5,159,0
144 Nov-11,,6,50,39,15,105,48,964,165,639,Data are confidential,0,18,16,2,36,6,147,3,59,0,1,13,10,4,27,9,25,4,131,0
145 Dec-11,,5,42,32,12,85,34,864,153,574,Data are confidential,0,16,16,2,34,5,120,3,56,0,0,11,9,4,24,8,24,2,113,0
146 Jan-12,,5,55,45,15,115,46,825,165,721,25,0,20,18,2,40,6,129,2,64,0,0,15,12,5,32,9,23,3,155,0
147 Feb-12,,6,48,37,12,97,34,658,135,592,19,0,18,15,2,34,4,110,2,52,0,0,12,10,4,27,7,18,3,124,0
148 Mar-12,,7,49,37,13,99,31,694,130,598,21,0,18,14,2,34,4,108,2,49,0,0,11,9,4,25,6,15,2,124,0
149 Apr-12,,5,60,43,17,120,38,803,149,724,24,0,22,16,3,41,5,122,2,58,0,0,15,11,5,32,7,20,3,153,0
150 May-12,,3,47,34,16,98,32,681,118,583,19,0,18,12,3,34,5,60,2,48,0,0,12,9,5,26,7,23,3,123,0
151 Jun-12,,3,42,30,17,90,31,668,119,570,19,0,16,11,4,32,5,84,2,49,0,0,10,7,5,22,7,30,2,120,0
152 Jul-12,,4,52,38,23,113,45,982,169,744,26,0,19,13,5,38,7,126,2,61,0,0,13,9,6,28,10,41,4,153,0
153 Aug-12,,5,41,30,17,88,34,892,145,600,21,0,14,10,3,28,5,112,2,52,0,0,13,8,5,26,8,45,3,129,0
154 Sep-12,,8,45,31,16,91,40,873,143,610,24,0,17,11,3,31,6,123,2,49,0,0,16,9,4,29,10,44,4,128,0
155 Oct-12,,9,60,43,19,122,58,1270,212,793,27,0,21,17,3,41,7,142,3,50,0,1,19,11,5,36,14,53,4,162,0
156 Nov-12,,7,48,36,15,100,49,912,147,672,21,0,16,14,2,33,6,119,2,27,0,1,13,10,4,28,11,41,3,133,0
157 Dec-12,,6,40,30,12,82,35,917,152,628,17,0,15,14,2,31,5,104,2,23,0,0,12,10,4,26,9,32,3,115,0
158 Jan-13,,7,52,41,15,108,48,937,182,762,25,0,20,18,2,40,6,134,2,29,0,1,15,13,5,33,10,31,4,155,0

```

- My data won't load ...
- ... because nobody bothered to use escape symbols.
- ... because ` is not a proper quotation symbol.
- ... because the maximum line length is exceeded.
- ... because there is a header row.
- ... because there is no header row.
- ... because the first line is the table-name.
- ... because some lines are empty.
- ... because it is encoded in CP-1252.
- ... because columns are shifted every ten rows.
- ... because a numeric column contains a string in line 590450.
- ... because some lines are two fields shorter.
- ... because Ümlauts are not supported.
- ... because someone added footnotes.
- ... because who uses § as a delimiter?
- ... because the file contains multiple tables.
- ... because tab and space are not the same thing.
- ... because someone added a comment in line 3.
- ... because - is not -.
- ... because it is split across multiple files.
- ... because headers are repeated every 80 lines.
- ... because the file ends mid-row.

Data Preparation: Tasks and Tools

- Data discovery
- Data validation
- Data structuring
- Data enrichment
- Data filtering
- Data cleaning

- And for data scientists
 - Feature selection
 - Feature extraction

| Categories | Available features | Data preparation tools | | | | | | |
|------------------|-------------------------------------|------------------------|--------|-----|-----|---------|--------|----------|
| | | Altair | Paxata | SAP | SAS | Tableau | Talend | Trifacta |
| Data discovery | Locate missing values (nulls) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Locate outliers | | ✓ | | ✓ | | | ✓ |
| | Search by pattern | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Sort data | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Data validation | Compare values (selection and join) | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ |
| | Check data range | ✓ | ✓ | ✓ | | ✓ | | |
| | Check permitted characters | | | | | | | |
| | Check column uniqueness | ✓ | ✓ | ✓ | | ✓ | | |
| Data structuring | Find type-mismatched data | | ✓ | ✓ | | ✓ | | |
| | Find data-mismatched datatypes | | ✓ | | | ✓ | | |
| | Change column data type | ✓ | ✓ | ✓ | ✓ | | | |
| | Delete column | ✓ | ✓ | ✓ | ✓ | | | |
| | Detect & change encoding | | | | | | | |
| | Pivot / unpivot | ✓ | ✓ | ✓ | | | | |
| | Rename column | ✓ | ✓ | ✓ | ✓ | | | |
| Data enrichment | Split column | ✓ | ✓ | ✓ | ✓ | | | |
| | Transform by example [13] | | | | | | | |
| | Assign semantic data type | | | | ✓ | ✓ | | |
| | Calculate column using expressions | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Discover & merge external data | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ |
| | Duplicate column | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ |
| | Generate primary key column | | | ✓ | | | | ✓ |
| | Join & union | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Merge columns | ✓ | | ✓ | | ✓ | ✓ | ✓ |
| | Normalize numeric values | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Data filtering | Delete/keep filtered rows | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Delete empty and invalid rows | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Extract value parts | ✓ | | | ✓ | | ✓ | ✓ |
| | Filter with regular expressions | | | | | | | ✓ |
| Data cleaning | Change date & time format | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Change letter case | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Change number format | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Deduplicate data | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ |
| | Delete by pattern | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ |
| | Edit & replace cell data | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Fill empty cells | ✓ | ✓ | | | | ✓ | ✓ |
| | Remove extra whitespace | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Remove diacritics | | | ✓ | | | | |
| | Standardize strings by pattern | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Standardize values in clusters | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |



Felix Naumann
Data Quality

Selected Data Preparation Projects – Bringing Order to Files

- Mondrian
 - Dissecting multi-table files
- ExracTable
 - Parsing visually delimited files
- Suragh and Tasheeh
 - Identifying ill-formed records
- Strudel
 - Classify cell-types
- AggreCol
 - Identify aggregation cells
- Pollock benchmark
 - Evaluate data ingestion ability



Felix Naumann
Data Quality

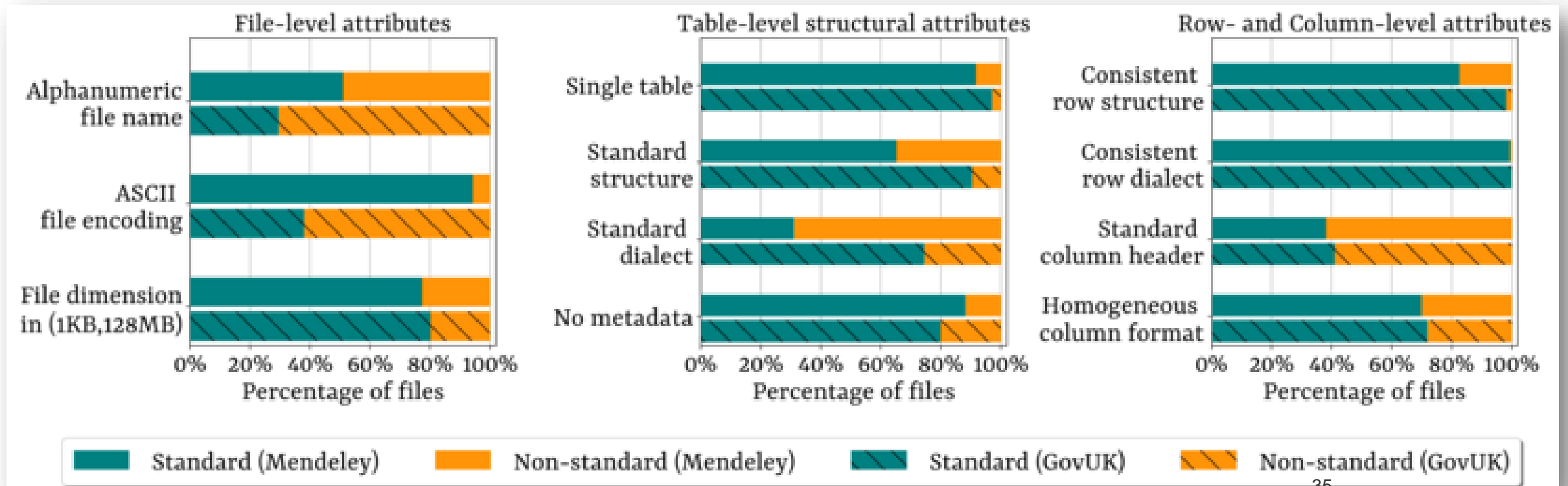
Pollock: Benchmarking the Ingestion Ability of Systems

```
Python 3.8.5 (default, Sep 3 2020, 21:29:08) [MSC v.1916 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import pandas as pd
>>> pd.read_csv("11-708-data-nlss-2009-1.csv")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "C:\Users\User\miniconda3\envs\pollution\lib\site-packages\pandas\io\parsers.py", line 686, in read_csv
    return _read(filepath_or_buffer, kwds)
  File "C:\Users\User\miniconda3\envs\pollution\lib\site-packages\pandas\io\parsers.py", line 458, in _read
    data = parser.read(nrows)
  File "C:\Users\User\miniconda3\envs\pollution\lib\site-packages\pandas\io\parsers.py", line 1196, in read
    ret = self._engine.read(nrows)
  File "C:\Users\User\miniconda3\envs\pollution\lib\site-packages\pandas\io\parsers.py", line 2155, in read
    data = self._reader.read(nrows)
  File "pandas\_libs\parsers.pyx", line 847, in pandas._libs.parsers.TextReader.read
  File "pandas\_libs\parsers.pyx", line 862, in pandas._libs.parsers.TextReader._read_low_memory
  File "pandas\_libs\parsers.pyx", line 918, in pandas._libs.parsers.TextReader._read_rows
  File "pandas\_libs\parsers.pyx", line 905, in pandas._libs.parsers.TextReader._tokenize_rows
  File "pandas\_libs\parsers.pyx", line 2042, in pandas._libs.parsers.raise_parser_error
pandas.errors.ParserError: Error tokenizing data. C error: Expected 25 fields in line 97, saw 27
```

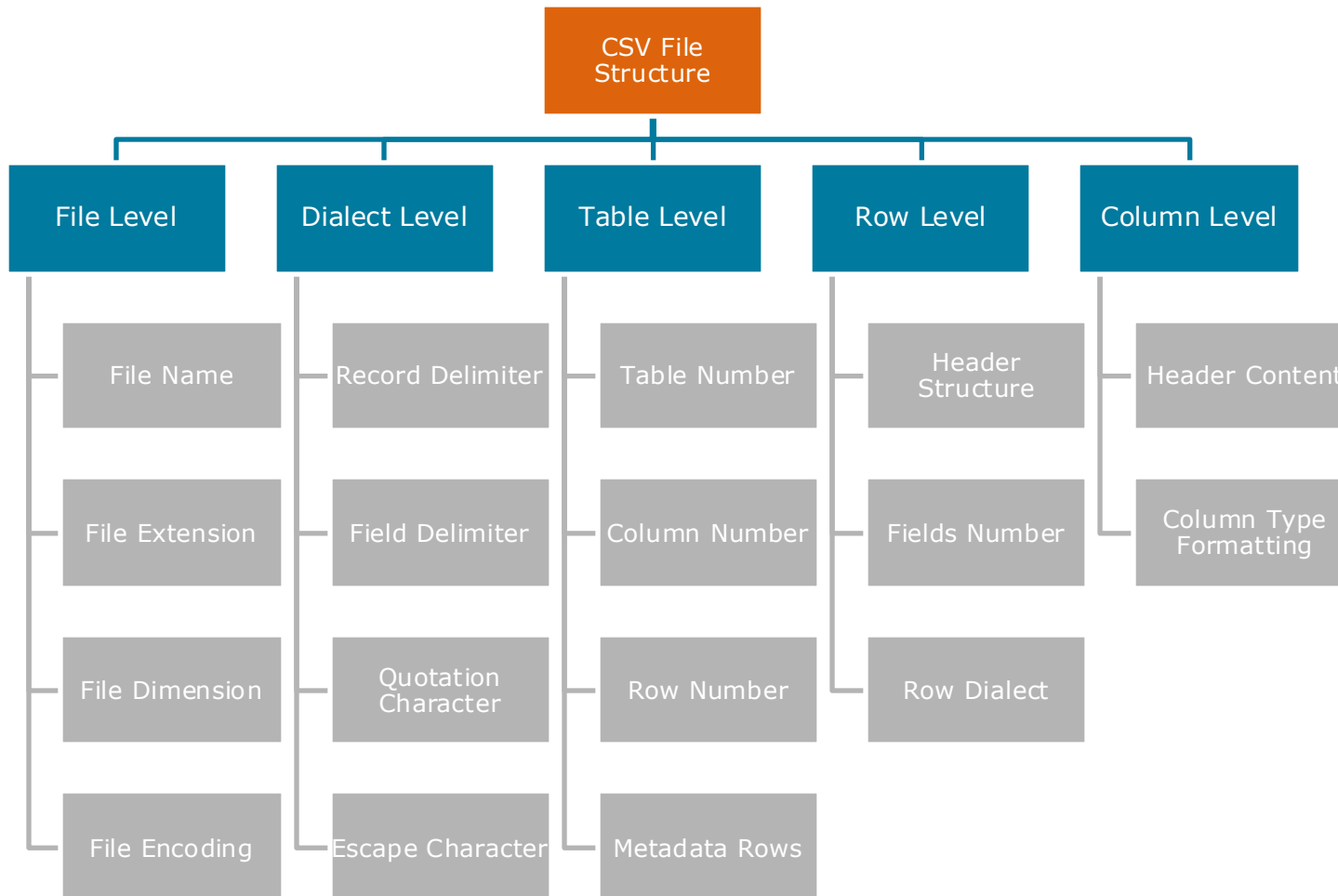
Felix Naumann
Data Quality

Pollock: Raw Data Survey

- Manual Annotation
 - 1,438 random files from GovUK
 - 2,274 random files from Mendeley



Pollock: Benchmark Dimensions and Results



| | Pollock score (2 289 +1) files | |
|------------------|-----------------------------------|-------------|
| | Simple | Weighted |
| CLEVERCSV 0.7.4 | 9.05 | 9.49 |
| CSVCOMMONS 1.9.0 | 6.63 | 9.29 |
| HYPOPARSR 0.1.0 | 3.73 | 4.41 |
| OPENCsv 5.6 | 6.62 | 7.80 |
| PANDAS 1.4.3 | 9.88 | 9.75 |
| PyCsv 3.10.5 | 9.71 | 9.47 |
| RCsv 4.2.1 | 7.78 | 6.76 |
| UNIVOCITY 2.9.1 | 9.35 | 7.97 |
| MARIADB 10.9.3 | 8.81 | 7.44 |
| MYSQL 8.0.31 | 8.88 | 7.45 |
| POSTGRESQL 15.0 | 0.14 | 7.33 |
| SQLITE 3.39.0 | 9.94 | 9.73 |
| CALC 7.3.6 | 9.75 | 7.52 |
| SPREADDESKTOP | 9.79 | 9.29 |
| SPREADWEB | 9.65 | 9.29 |
| DATAVIZ | 4.93 | 5.51 |

Agenda

1. Data and Information Quality Research
2. Data Preparation
- 3. Data Quality and AI Systems**
 - With Hazar Harmouch, Sedir Mohammed et al.
4. Data Quality Assessment



Empirical Measurement of the Effects of Poor Data Quality on ML Results: Measurement Dimensions

Pollutions

- Consistent representation
- Completeness
- Feature accuracy
- Target accuracy
- Uniqueness
- Target balance

Scenarios

- Pollute only training data
- Pollute only test data
- Pollute training and test data

Runs

- 5 runs, average

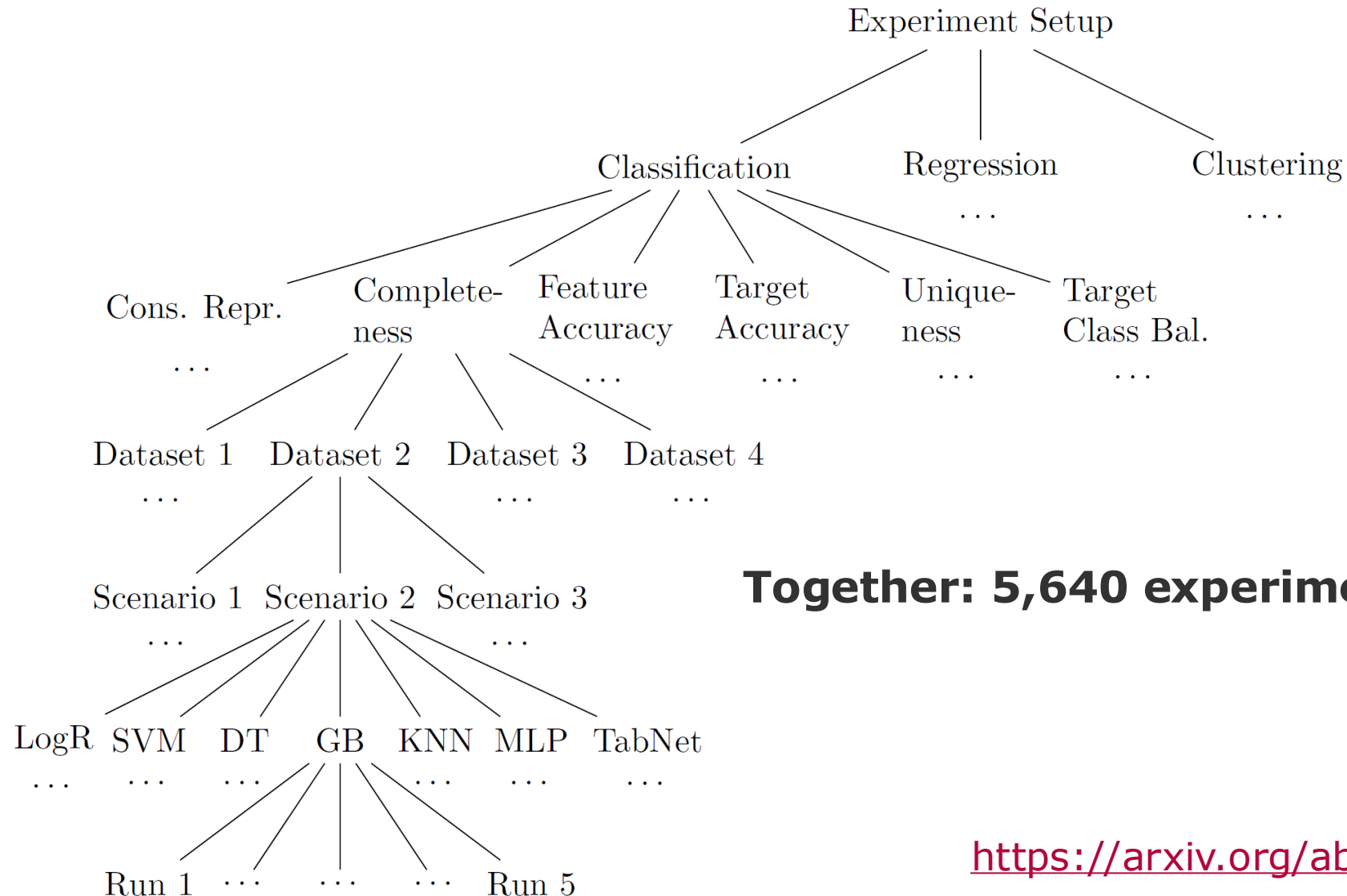
Tasks and algorithms

- Classification
 - LogR, SVM, DT, GB, KNN, MLP
- Clustering
 - GM, k-Means, k-Prototypes, AC, OPTICS
- Regression
 - LR, RR, DT, RF, GB, MLP, TabNet

Datasets

- TelcoChurn, GermanCredit, Contraceptive, COVID
- Houses, IMDB, Cars
- Bank, Covertypes, Letter

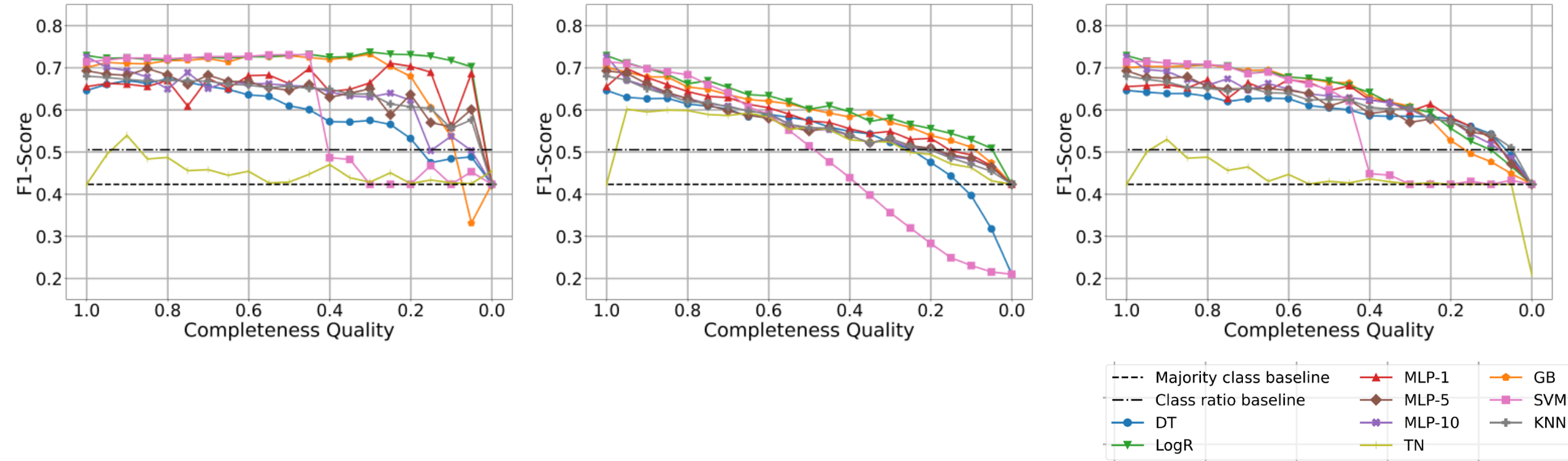
Empirical Measurement of the Effects of Poor Data Quality on ML Results: Measurement Dimensions



Together: 5,640 experiments

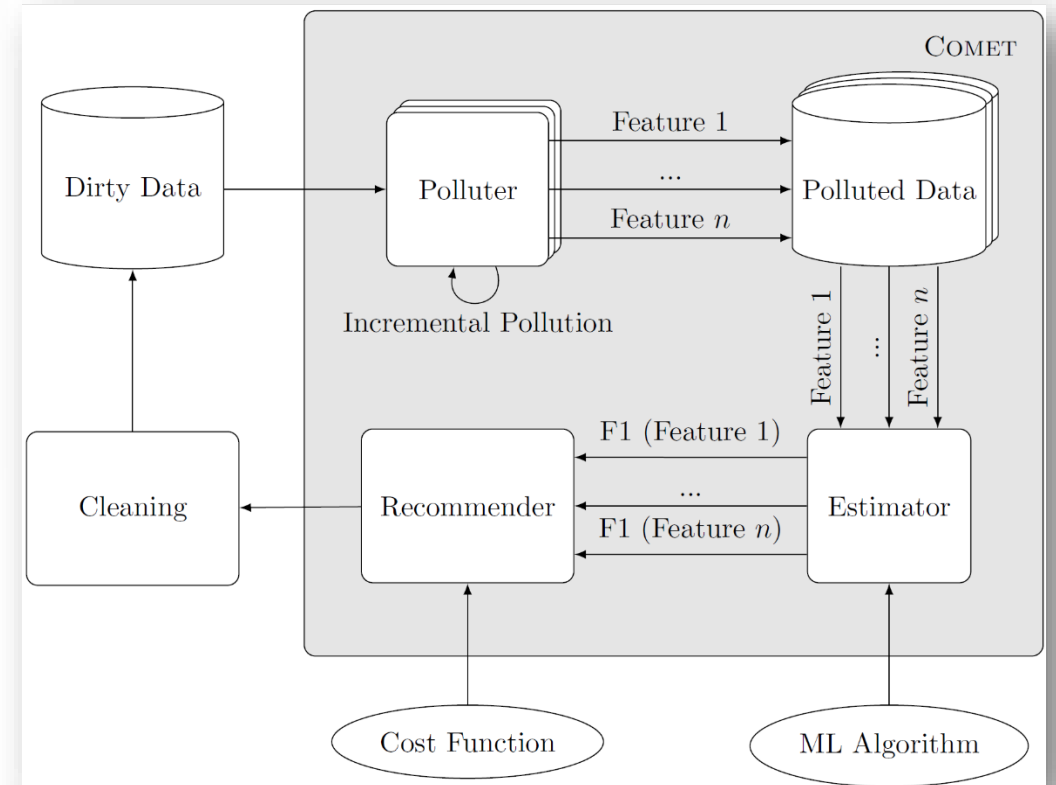
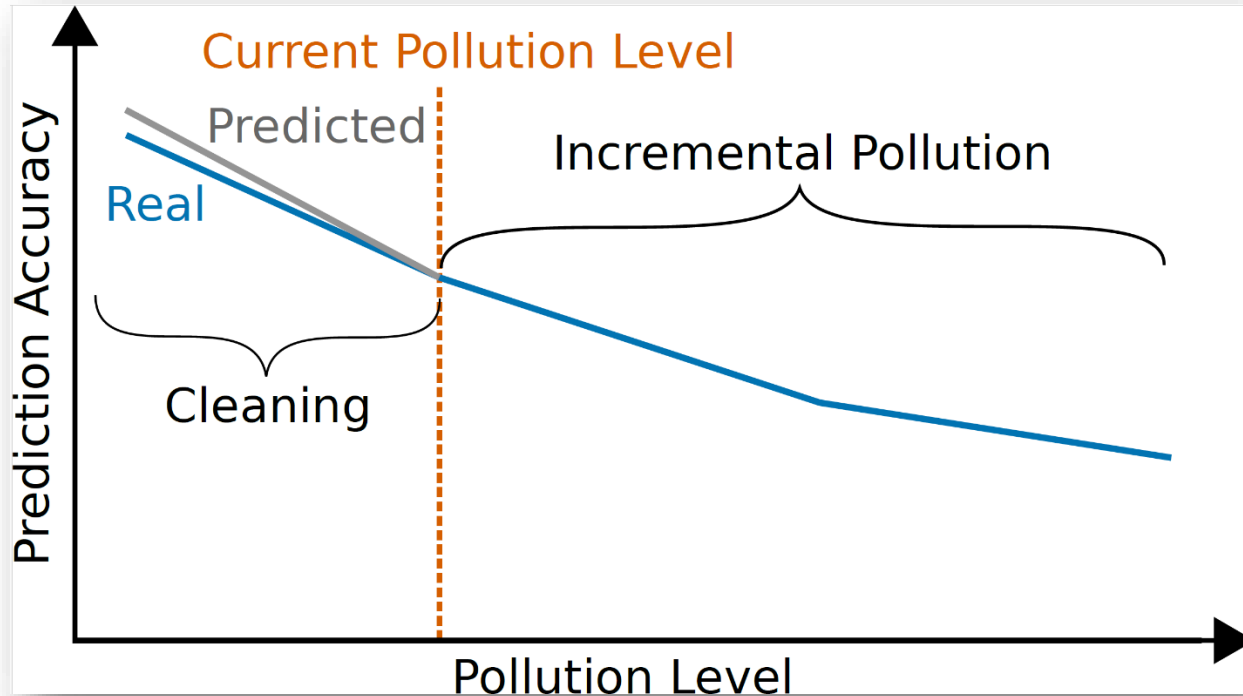
Felix Naumann
Data Quality

Example Results



Average F1-Score for Classification of the Telco-Churn dataset

So what? Recommend Data Cleaning Steps



Agenda

1. Data and Information Quality Research
2. Data Preparation
3. Data Quality and AI Systems
4. **Data Quality Assessment**
 - With Hazar Harmouch, Lisa Ehrlinger, Sedir Mohammed and Divesh Srivastava

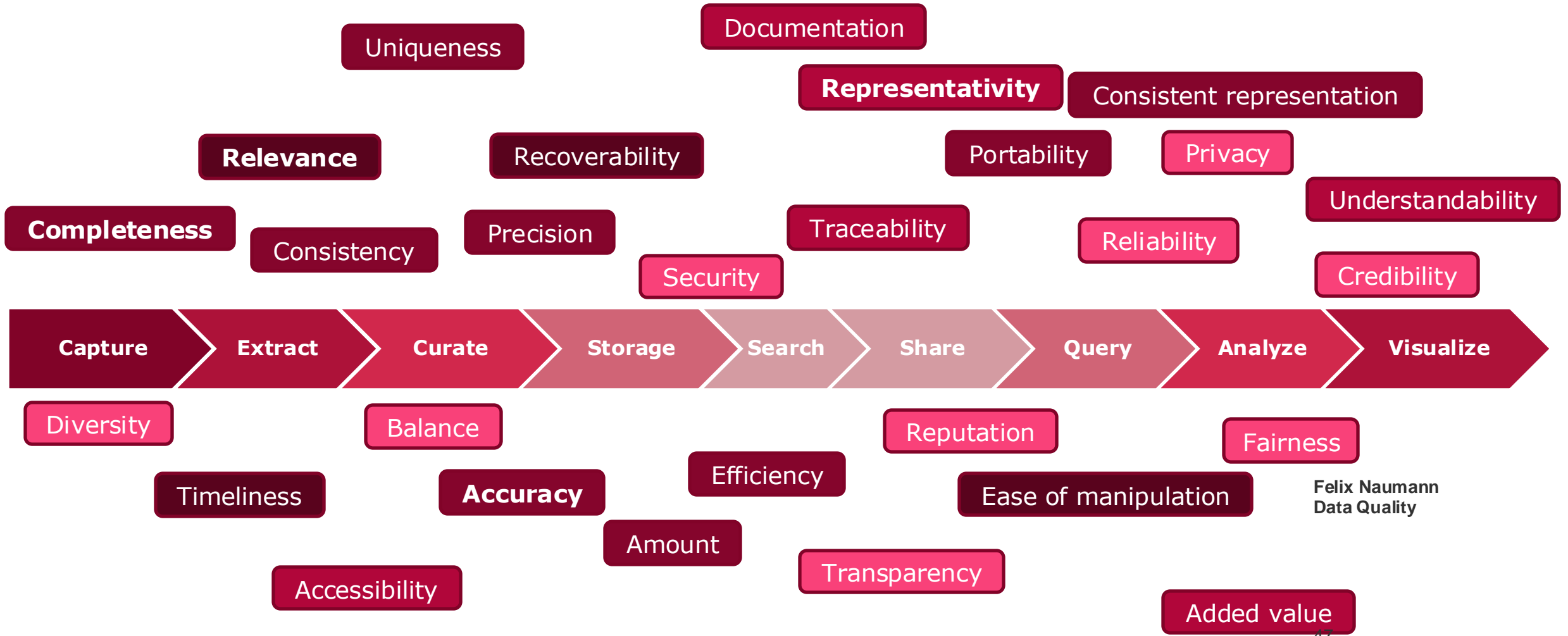


European AI Act Article 10 (3): Data and Data Governance

- **High-quality data** and access to high-quality data plays a vital role in providing structure and in ensuring the performance of many AI systems, especially when techniques involving the training of models are used, with a view to ensure that the high-risk AI system performs as intended and safely and it does not become a source of discrimination prohibited by Union law.
- High-quality data sets for training, validation and testing require the implementation of appropriate **data governance and management** practices.
- Data sets for training, validation and testing, including the labels, should be **relevant**, sufficiently **representative**, and to the best extent possible **free of errors** and **complete** in view of the intended purpose of the system.
- The data sets should also have the **appropriate statistical properties**, including as regards the persons or groups of persons in relation to whom the high-risk AI system is intended to be used, with specific attention to the mitigation of possible biases in the data sets [...].



Data Quality along the AI Pipeline



Felix Naumann
Data Quality

DQ Assessment in the Year 2000

■ [c5]    Felix Naumann, Claudia Rolker:

Assessment Methods for Information Quality Criteria. IQ 2000: 148-162

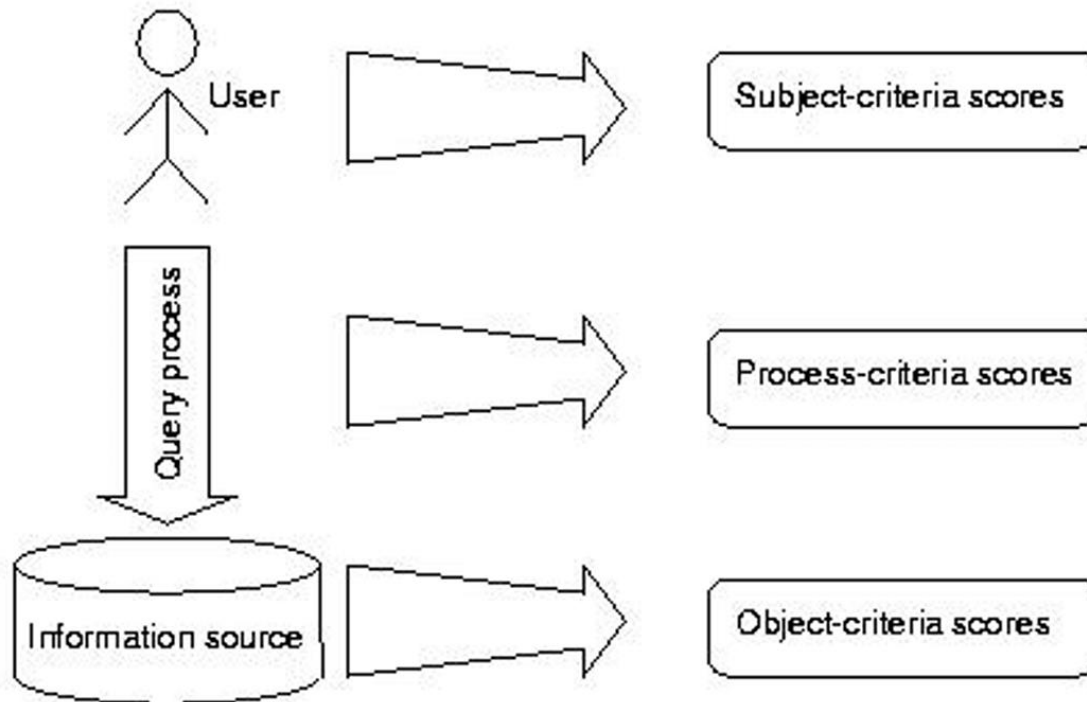
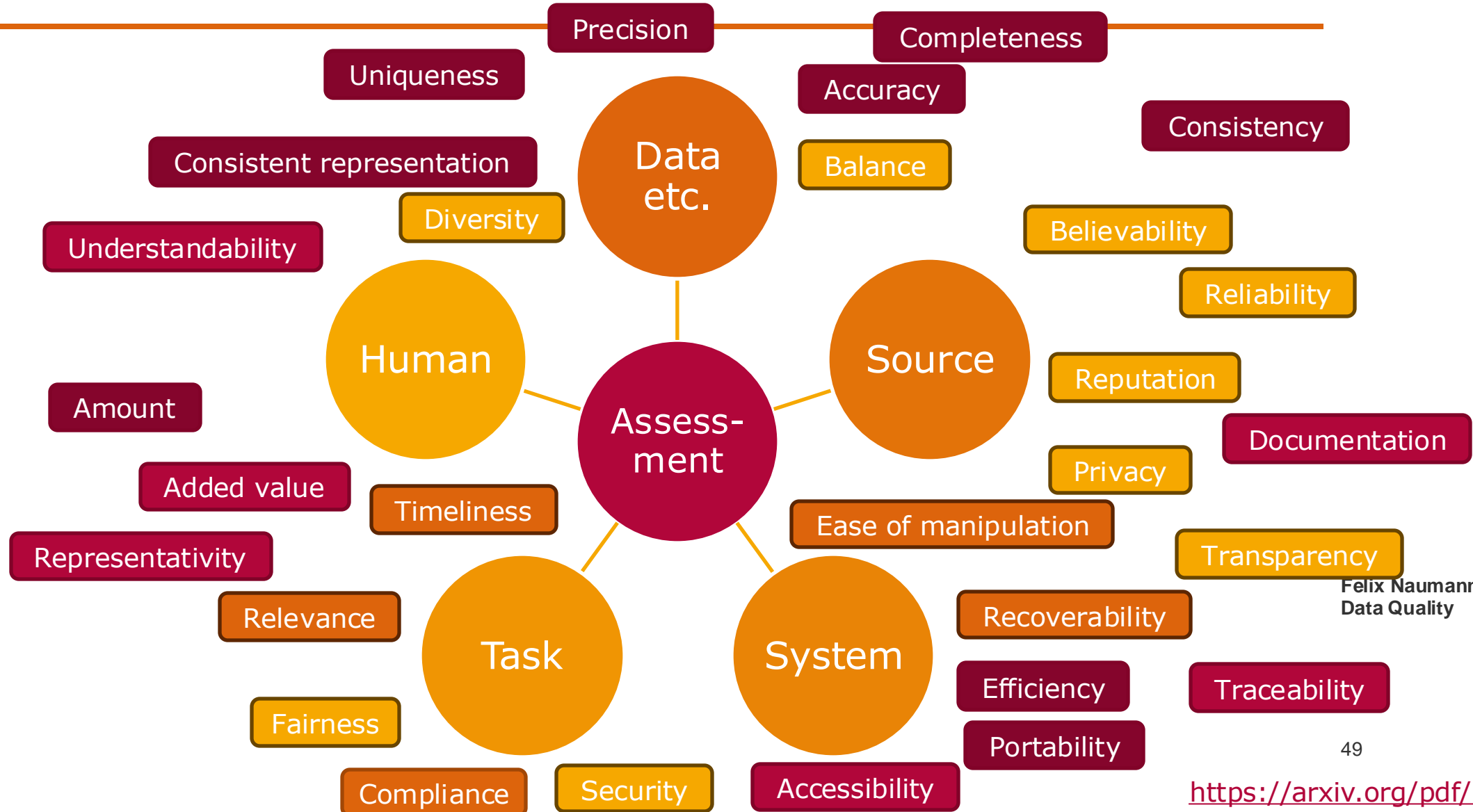


Figure 1: Three sources of IQ criterion scores

| Assessment Class | IQ Criterion | Assessment Method |
|------------------|---------------------------|--------------------------------|
| Subject Criteria | Believability | User experience |
| | Concise representation | User sampling |
| | Interpretability | User sampling |
| | Relevancy | Continuous user assessment |
| | Reputation | User experience |
| | Understandability | User sampling |
| | Value-Added | Continuous user assessment |
| Object Criteria | Completeness | Parsing, sampling |
| | Customer Support | Parsing, contract |
| | Documentation | Parsing |
| | Objectivity | Expert input |
| | Price | Contract |
| | Reliability | Continuous assessment |
| | Security | Parsing |
| | Timeliness | Parsing |
| Process Criteria | Verifiability | Expert input |
| | Accuracy | Sampling, cleansing techniques |
| | Amount of data | Continuous assessment |
| | Availability | Continuous assessment |
| | Consistent representation | Parsing |
| | Latency | Continuous assessment |
| | Response time | Continuous assessment |

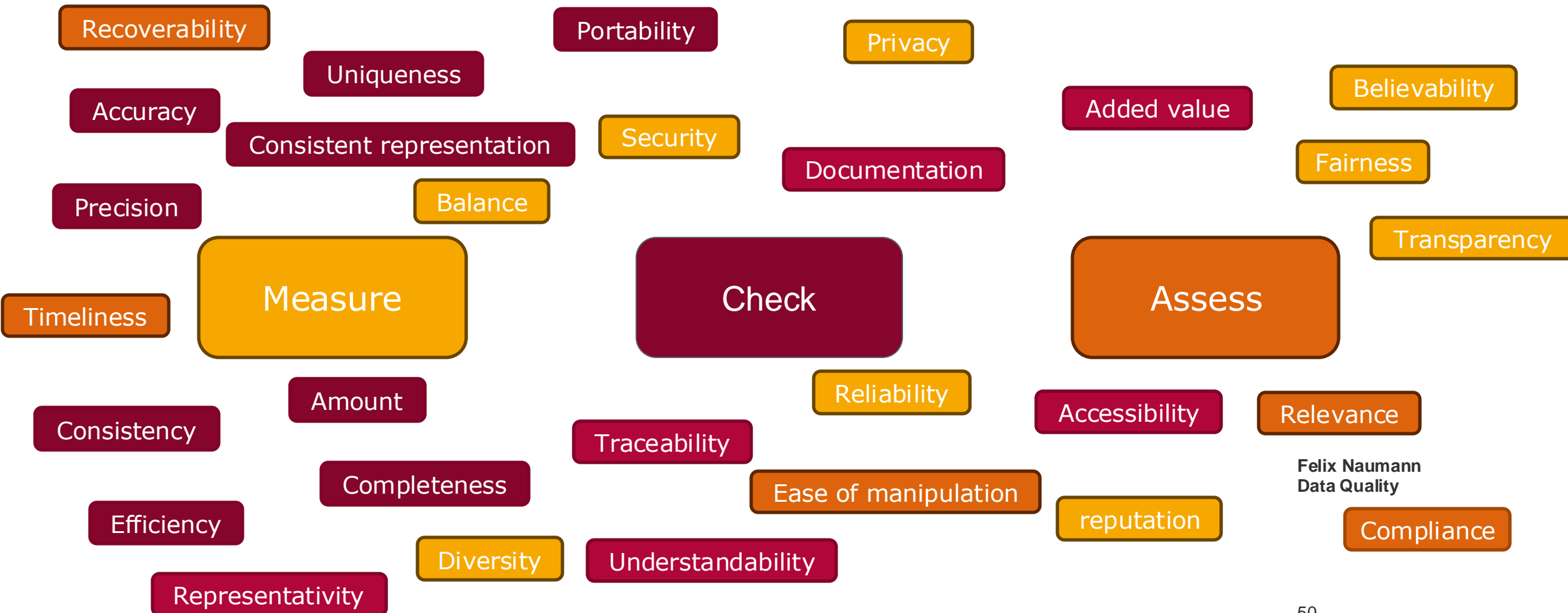
Table 2: Classification of IQ Metadata Criteria

Ingredients for DQ Assessment: Five Facets



Felix Naumann
Data Quality

Assessing Data Quality



Felix Naumann
Data Quality

Assessment Examples

Completeness

- Values vs. rows vs. columns
- Absolute (counts) or relative (percentages)
- Relative to what? External data needed
- Semantically challenging

Representativity

- vs. balance vs. diversity
- Presence of every value combination
 - Existing values vs. all values
 - Computationally challenging
- Distribution similar to real-world distribution: external data needed

Free-of-errors / Correctness

- Error detection
 - Count at value or row-level
- Business rules
 - Patterns, dependencies, data-types
- Outlier detection
- Validation with external data

Relevance

- ...

Understandability

- ...

Further Challenges for DQ Assessment

■ Ambiguity

- Many attempts to compile and define DQ dimensions
- Definitions of the dimensions inherently ambiguous

■ Explainability

- Assessment results explainable to consumers
- Results traceable to their root cause, to improve quality

■ Efficiency

- Assessment effort and time should be low

■ Compliance

- Fulfill organizational data governance processes
- Comply to a legal framework, e.g., GDPR or the AI Act

■ Scoring

- Aggregate and normalize assessment results to some numeric scale.
- Allows comparison across datasets and across time

■ Adequacy

- Is the data of sufficient quality or adequate for the task at hand?

Summary

- Data and Information Quality Research
- Data Preparation
- Data Quality and AI Systems
- Data Quality Assessment

