



Data Quality in the Age of AI

Keynote at Guide-AI workshop @SIGMOD 2024

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

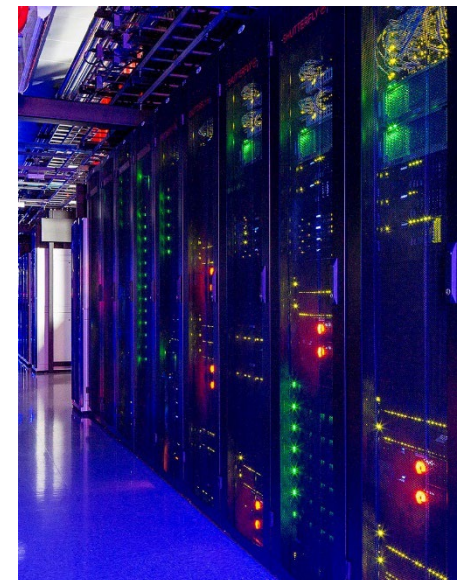
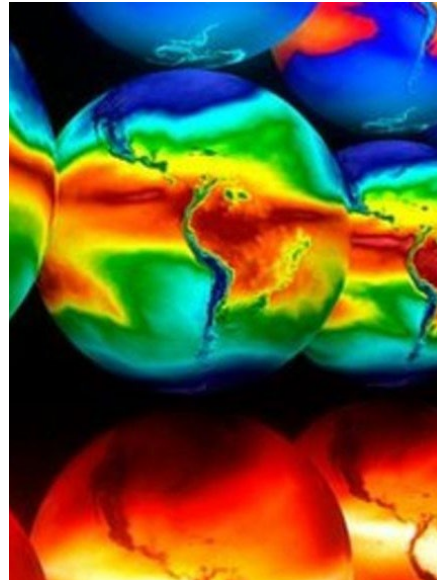
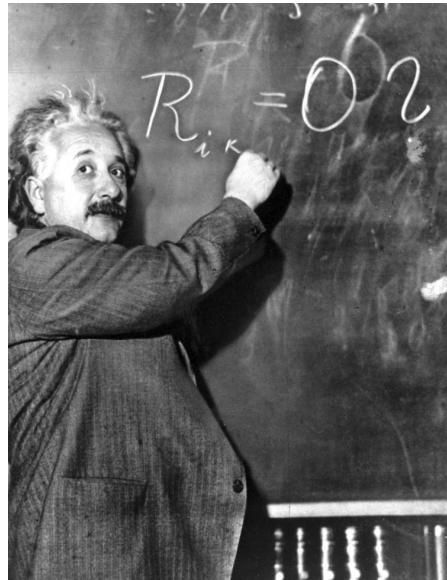
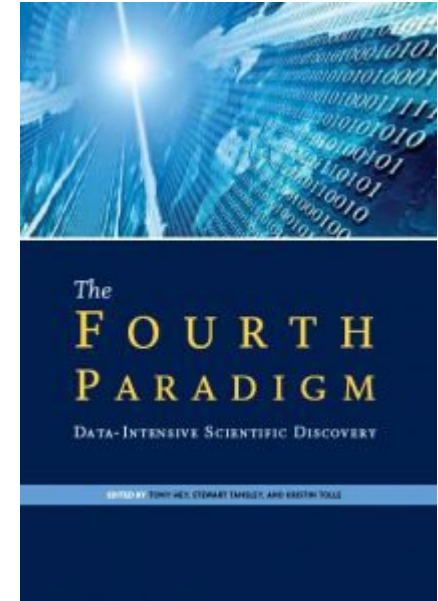
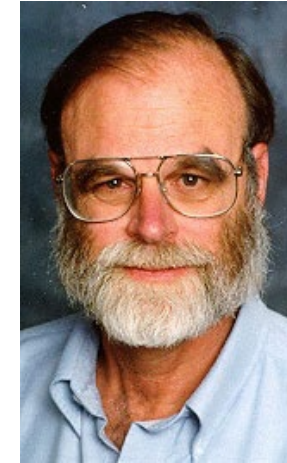
June 14, 2024

Felix Naumann

The Fourth Paradigm of Science

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

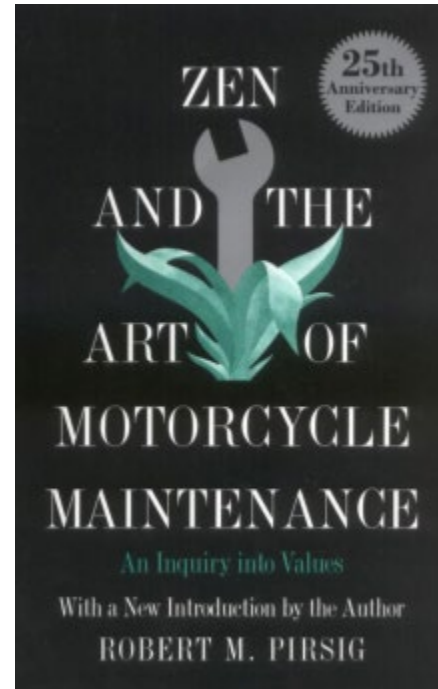
*We have to do better producing tools to **support the whole research cycle** - from data capture and data curation to data analysis and data visualization. Jim Gray*



Felix Naumann
Data Quality

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

Robert Pirsig





Felix Naumann
Data Quality

Data Errors for Database Researchers



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-	I-
- 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

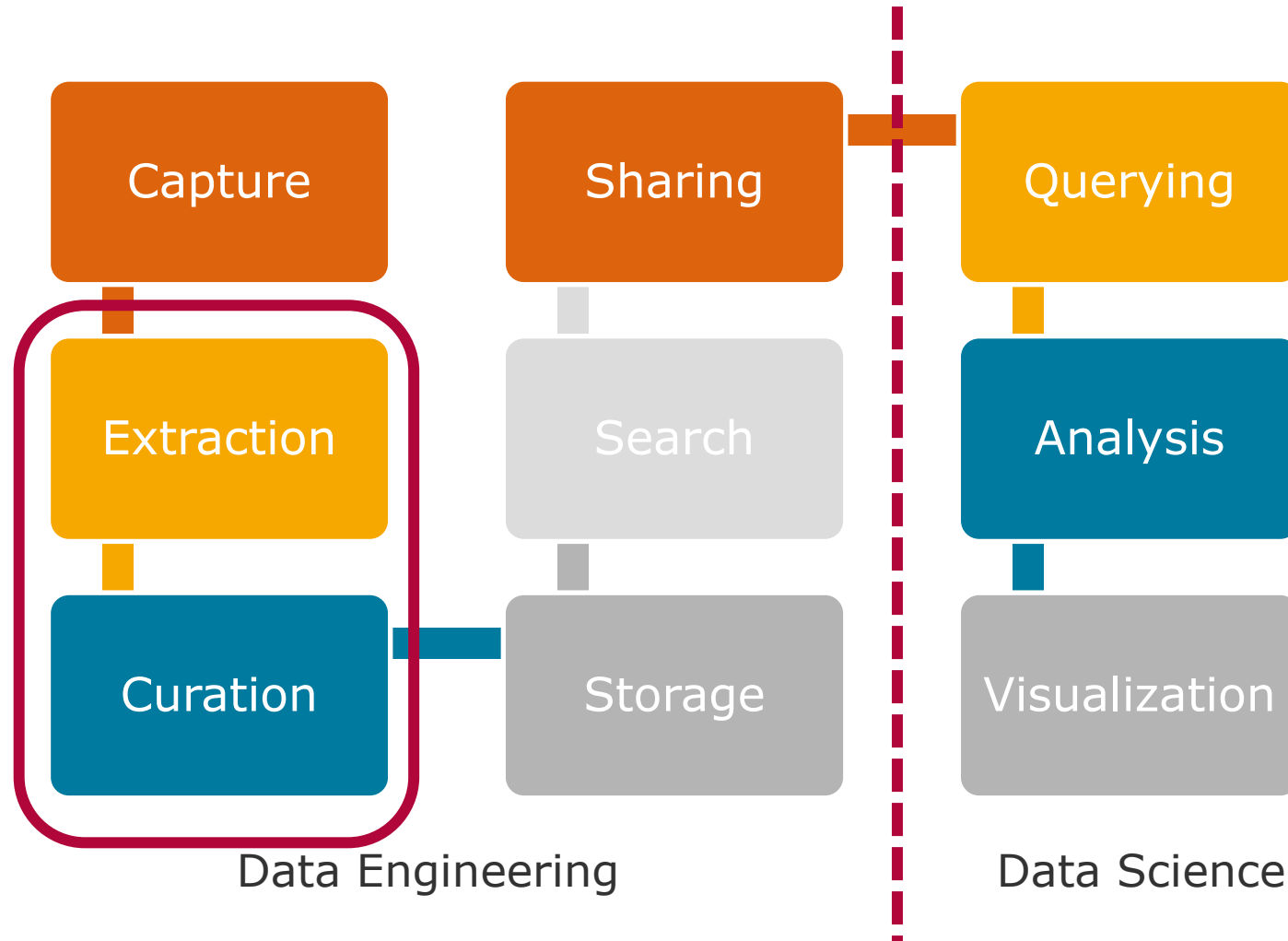
Incorrect addresses
stall nearly 1,500
Tennessee refunds

By **BONNA de la CRUZ**
Staff Writer

Now that Tilcia L. Meniffee 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

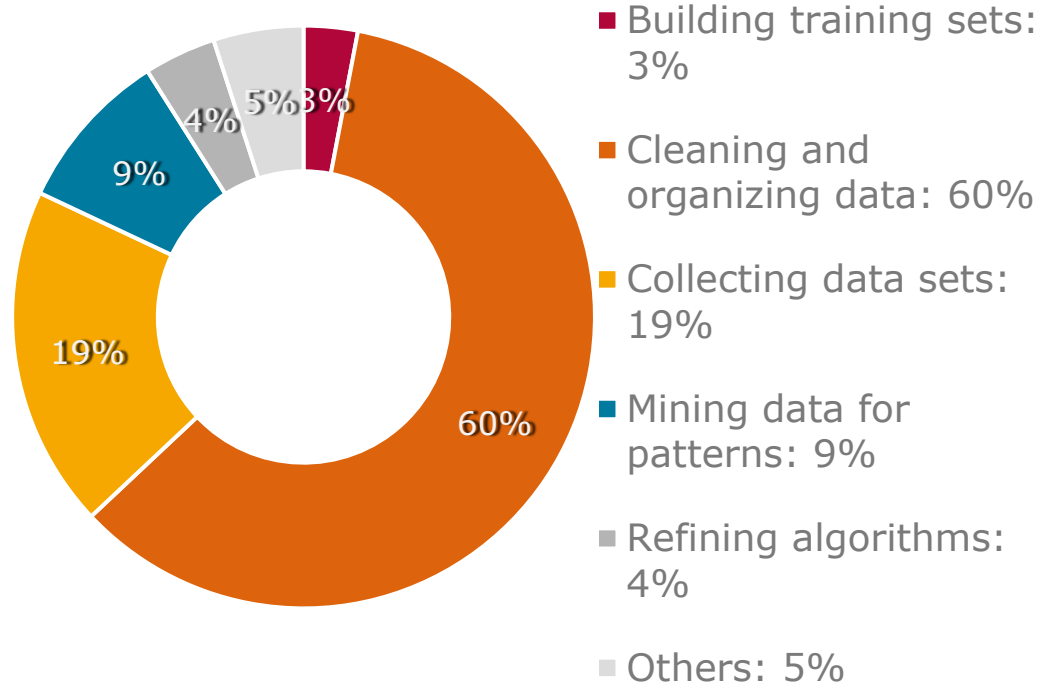
Data Science Pipeline



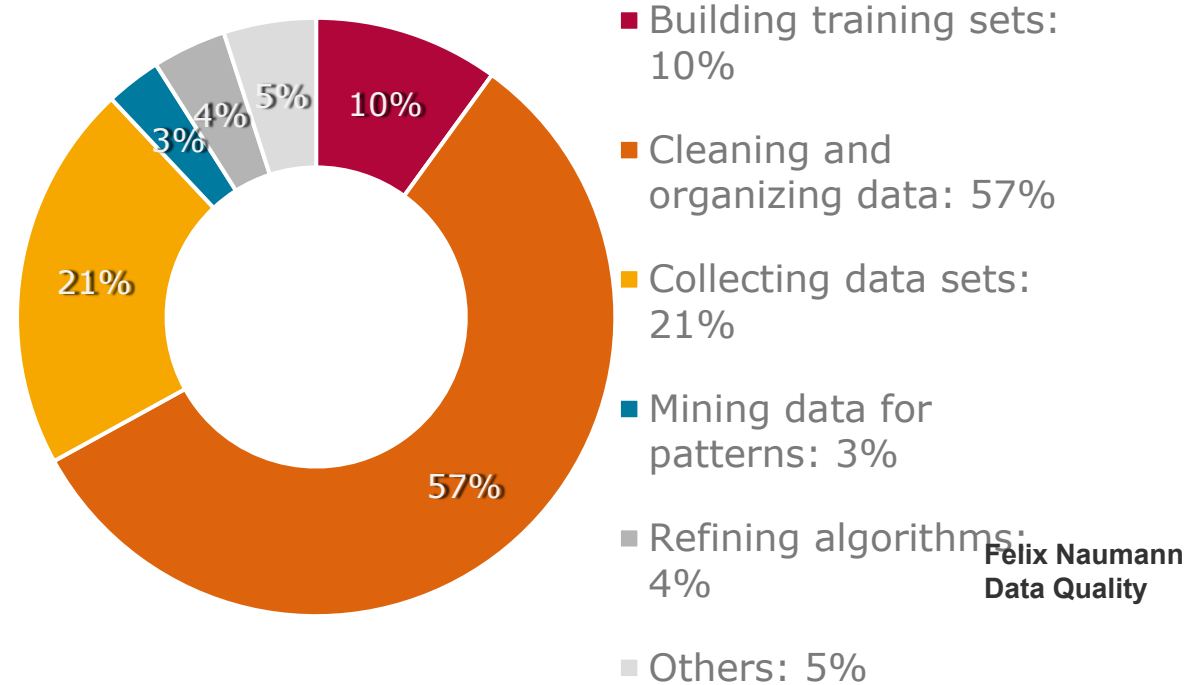
Felix Naumann
Data Quality

Data preparation in reality

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

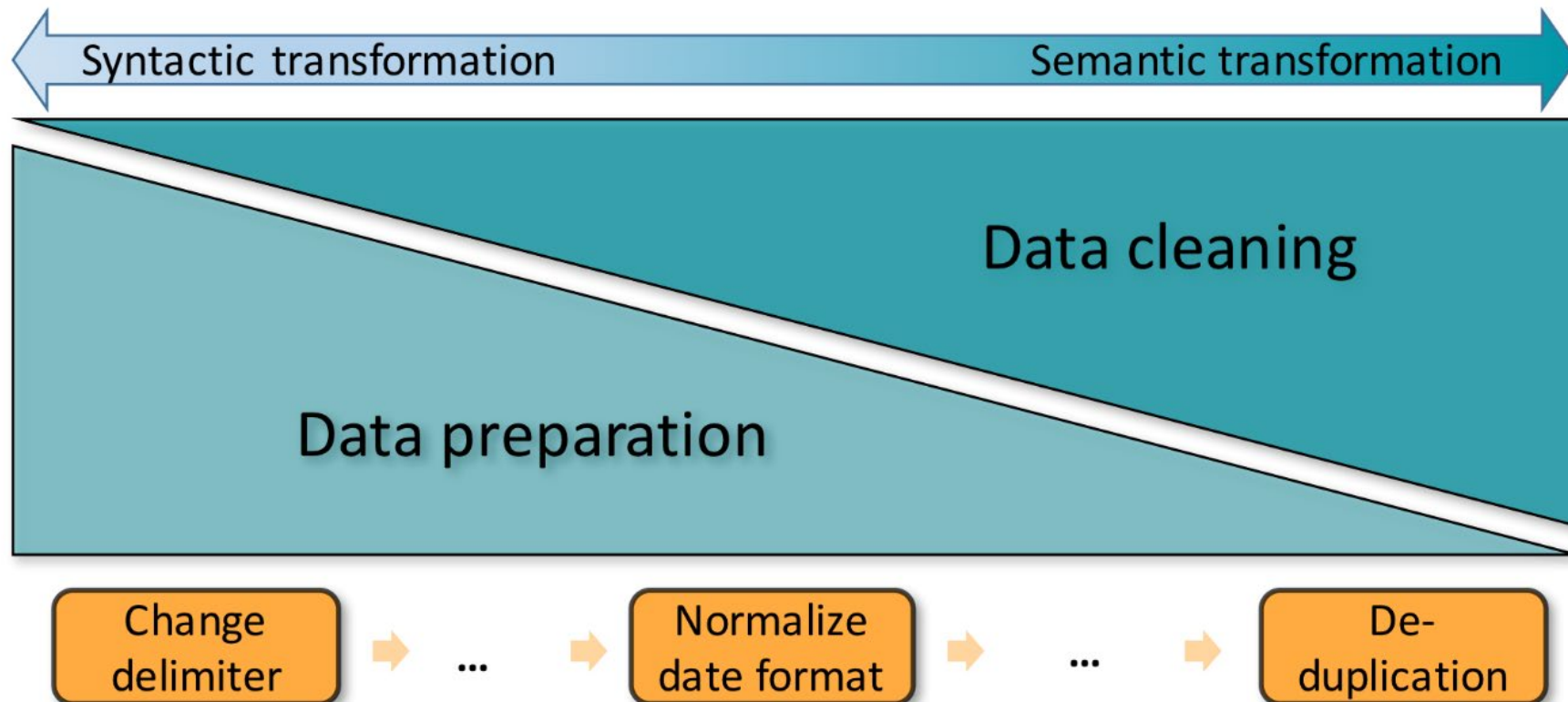


What is the **least enjoyable** part of data science?



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



Zooming into Information Quality



Fitness for use

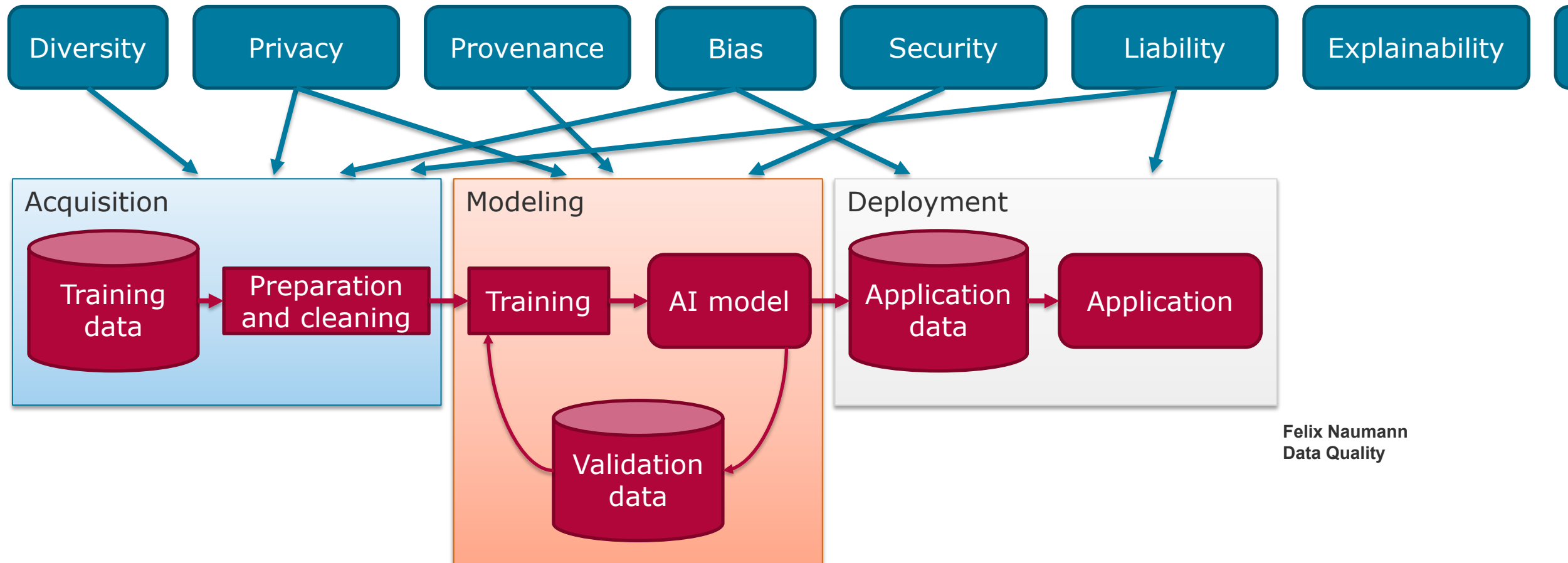
Accuracy, Objectivity, Believability, Reputation, Accessibility, Security, Relevance, Value-Added, Timeliness, Completeness, Amount of Data, Interpretability, Understandability, Consistency, Concise Representation

179 Dimensions

Ability to be Joined With	Ability to Download	Ability to Identify Errors	Ability to Upload	Extensibility	Extern	Finalization	Flawlessness
Acceptability	Access by Comparison	Accessibility	Accuracy	Flexibility	Form of Presentation	Format	Irregularity
Adaptability	Adequate Detail	Adequate Volume	Ample Information	Friendliness	Generality	Habit	Historical
Age	Aggregability	Availability	Amount of Data	Importance	Inconsistencies	Integration	Occurrence
Auditable	Authority	Believability	Be-levability	Interactivity	Interesting	Level of Abstraction	Imagery
Breadth of Data	Browser	Clarity	Clarity	Localized	Logically Connected	Manageability	Level of Standardization
Clarity of Origin	Case Data	Conciseness	Conciseness	Measurability	Medium	Meets Requirements	Mapping
Comprehensibility	Completeness	Consistency	Consistency	Modularity	Narrowly Defined	No last information	Minimality
Competitive Edge	Confidence	Continuity	Continuity	Novelty	Objectivity	Optimality	Normality
Consistency	Conformity	Cost	Cost	Origin	Originality	Partitionability	Orderliness
Convenience	Correctness	Critical	Critical	Pedagogy	Personalized	Pertinence	Past Experience
Cost of Accuracy	Cost of Collection	Customizability	Customizability	Practicality	Practicality	Proprietary Nature	Purpose
Current	Current	Efficiency	Efficiency	Quantity	Rationality	Redundancy	Regularity of Form
Data Derived	Definability	Depth of Data	Depth of Data	Reliability	Reliability	Repetitive	Reproducibility
Data	Detailed Source	Distinguishability	Distinguishability	Relevance	Regulation of Graphics	Responsibility	Reliability
Dynamic	Ease of Access	Updated Files	Updated Files	Revisability	Revisability	Rigidity	Robustness
Ease of Data Exchange	Ease of Maintenance	Ease of Retrieval	Ease of Retrieval	Scope of Info	Secrecy	Security	Self-Correcting
Ease of Update	Ease of Use	Ergonomic	Ergonomic	Semantic	Self-Correcting	Size	Source
Efficiency	Efficiency	Expense	Expense	Interpretation	Speed	Stability	Storage
Error-Free	Expendability	Expendability	Expendability	Specificity	Speed	Stability	Storage
				Synchronization	Time-Independence	Timeliness	Traceable
				Translatable	Transparency	Unambiguity	Unbiased
				Understandable	Uniqueness	Unorganized	Up-to-Date
				Usable	Usefulness	Use Priority	Valid
				Value	Variability	Variety	Verifiable
				Votability	Well-Documented	Well-Preserved	

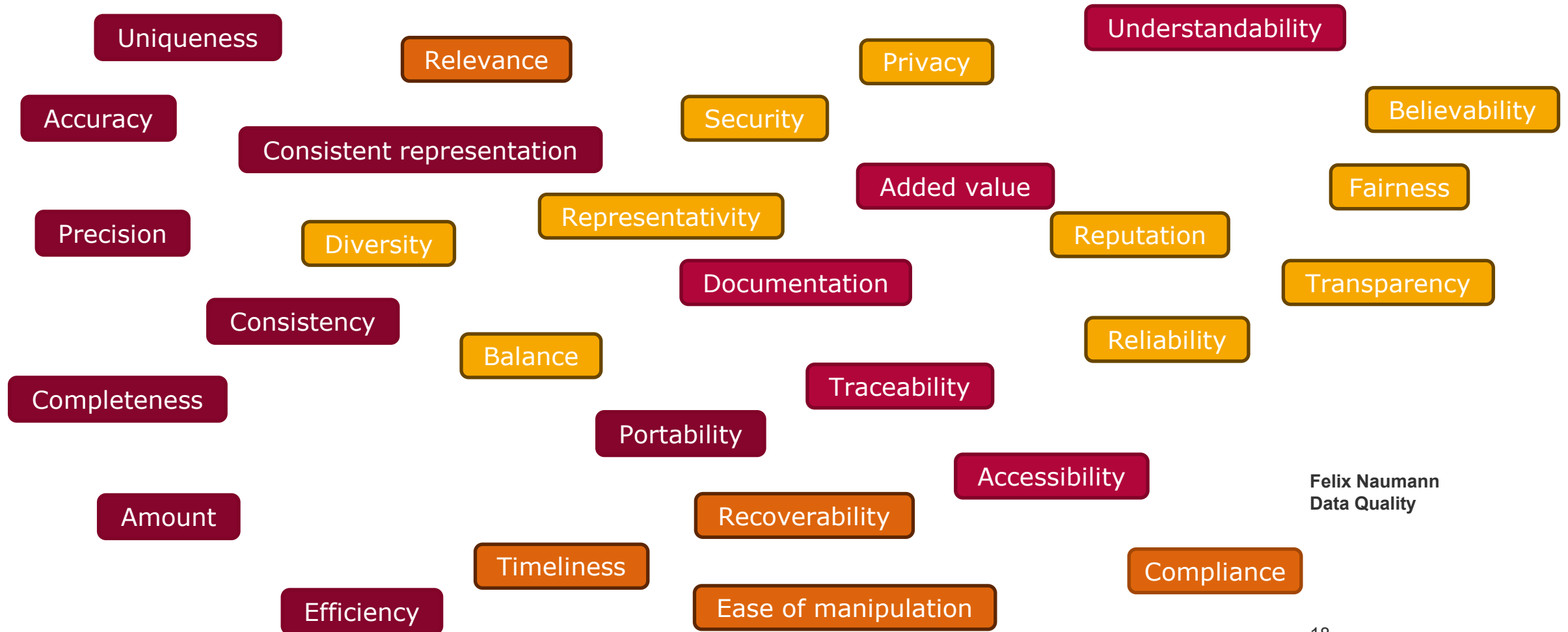
Felix Naumann
Data Quality

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 Aalten (-Strannood)	F	18	168		Netherlands	NED	1932 Summer	1932	Summer	Los Angeles	Athletics	Athletics Women's 100 metres	NA
28	Cornelia 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,,4,47,35,17,99,32,1055,185,578,16,0,18,16,2,36,5,149,2,47,0,0,16,11,5,32,10,43,5,115,1
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,* ,3,51,41,17,109,33,799,143,654,19,0,20,18,2,39,5,125,2,52,0,0,14,13,6,33,8,35,5,138,1
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,,3,43,33,20,97,26,682,132,531,14,0,18,13,5,36,4,55,1,39,0,0,11,10,8,29,6,27,5,115,1
128 Jul-10,* ,5,55,40,26,121,36,1075,182,662,Data are confidential,0,21,16,6,43,5,114,2,51,0,0,11,10,10,31,8,35,5,144,1
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,0,19,13,4,36,5,148,2,46,0,0,16,10,5,31,7,27,5,121,1
131 Oct-10,* ,9,63,44,22,129,49,1191,195,728,Data are confidential,0,24,19,4,47,6,197,3,57,0,1,22,13,6,41,10,29,7,157,1
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,0,21,20,3,43,5,167,3,54,0,0,14,11,6,31,8,20,4,143,1
134 Jan-11,* ,6,60,48,18,126,52,856,190,690,Data are confidential,0,22,20,3,45,6,148,2,52,0,1,16,15,7,38,10,19,4,157,1
135 Feb-11,,7,47,39,15,101,37,699,156,592,Data are confidential,0,19,16,2,37,4,115,2,48,0,0,14,12,5,32,8,13,2,123,1
136 Mar-11,,8,51,38,16,105,34,678,137,587,Data are confidential,0,20,16,2,37,4,115,2,49,0,0,13,11,5,29,6,12,2,122,1
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,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,8,30,4,127,0
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
159 Feb-13,,8,46,36,13,84,36,768,154,634,20,0,18,15,1,34,5,110,2,35,0,0,12,10,5,27,7,18,3,133,0
160 Mar-13,,9,53,42,19,127,37,827,167,683,Data are confidential,0,23,18,4,45,5,118,2,60,0,0,15,12,5,32,7,15,3,143,0
161 Apr-13,* ,7,62,46,19,127,37,827,167,683,Data are confidential,0,20,16,2,37,4,115,2,48,0,0,14,12,5,32,8,13,2,123,1
162 May-13,,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
163 Jun-13,,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
164 Jul-13,* ,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
165 Aug-13,,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
166 Sep-13,* ,7,51,36,17,104,40,992,153,621,Data are confidential,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,8,30,4,127,0
167 Oct-13,* ,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
168 Nov-13,,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
169 Dec-13,,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
170 Jan-14,* ,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
171 Feb-14,,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
172 Mar-14,,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
173 Apr-14,* ,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
174 May-14,,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
175 Jun-14,,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
176 Jul-14,* ,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
177 Aug-14,,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
178 Sep-14,,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
179 Oct-14,* ,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
180 Nov-14,,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
181 Dec-14,,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
182 Jan-15,* ,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
183 Feb-15,,8,46,36,13,84,36,768,154,634,20,0,18,15,1,34,5,110,2,35,0,0,12,10,5,27,7,18,3,133,0
184 Mar-15,* ,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
185 Apr-15,,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
186 May-15,,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
187 Jun-15,* ,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
188 Jul-15,,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
189 Aug-15,* ,7,51,36,17,104,40,992,153,621,Data are confidential,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,8,30,4,127,0
190 Sep-15,* ,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
191 Oct-15,,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
192 Nov-15,,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
193 Dec-15,* ,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
194 Jan-16,,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
195 Feb-16,,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
196 Mar-16,* ,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
197 Apr-16,,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
198 May-16,,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
199 Jun-16,* ,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
200 Jul-16,,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
201 Aug-16,,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
202 Sep-16,* ,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
203 Oct-16,,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
204 Nov-16,,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
205 Dec-16,* ,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
206 Jan-17,,8,46,36,13,84,36,768,154,634,20,0,18,15,1,34,5,110,2,35,0,0,12,10,5,27,7,18,3,133,0
207 Feb-17,* ,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
208 Mar-17,,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
209 Apr-17,,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
210 May-17,* ,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
211 Jun-17,,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
212 Jul-17,* ,7,51,36,17,104,40,992,153,621,Data are confidential,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,8,30,4,127,0
213 Aug-17,* ,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
214 Sep-17,,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
215 Oct-17,,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
216 Nov-17,* ,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
217 Dec-17,,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
218 Jan-18,,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
219 Feb-18,* ,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
220 Mar-18,,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
221 Apr-18,,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
222 May-18,* ,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
223 Jun-18,,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
224 Jul-18,,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
225 Aug-18,* ,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
226 Sep-18,,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
227 Oct-18,,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
228 Nov-18,* ,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
229 Dec-18,,8,46,36,13,84,36,768,154,634,20,0,18,15,1,34,5,110,2,35,0,0,12,10,5,27,7,18,3,133,0
230 Jan-19,* ,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
231 Feb-19,,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
232 Mar-19,,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
233 Apr-19,* ,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
234 May-19,,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
235 Jun-19,* ,7,51,36,17,104,40,992,153,621,Data are confidential,0,18,14,2,35,5,144,3,49,0,1,17,9,4,30,8,30,4,127,0
236 Jul-19,* ,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
237 Aug-19,,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
238 Sep-19,,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
239 Oct-19,* ,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
240 Nov-19,,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
241 Dec-19,,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
242 Jan-20,* ,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
243 Feb-20,,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
244 Mar-20,,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
245 Apr-20,* ,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
246 May-20,,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
247 Jun-20,,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
248 Jul-20,* ,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
249 Aug-20,,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
250 Sep-20,,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
251 Oct-20,* ,7,52,41,15,108,48,937,182,762
```

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	✓	✓	✓				
Data enrichment	Rename column	✓	✓	✓	✓			
	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	✓		✓		✓	✓	✓
Data filtering	Normalize numeric values	✓	✓	✓	✓	✓	✓	✓
	Delete/keep filtered rows	✓	✓	✓	✓	✓	✓	✓
	Delete empty and invalid rows	✓	✓	✓	✓	✓	✓	✓
	Extract value parts	✓			✓		✓	✓
Data cleaning	Filter with regular expressions							✓
	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



Felix Naumann
Data Quality

Mondrian: Multitable Spreadsheets

ECONOMIC CALCULATOR FOR FIRM VERSUS NON-FIRM PURCHASE			
Current Hour = 14	Next Hour = 15	Local Avail. 529	PV
Enter Local Generation Avail: 529	Projected Control Area Load: 840	Gen. 397	IID 80
Enter Remote Generation: 581	PNM Contingent: 41	UnLoaded 132	real/time 80
Firm Purchases into EPE: 25	TNP Firm: 25		net/pre -75
Non-Firm Purchases into EPE: 0	IID Firm + Contingent: 150		Tep/exc 300
SPS Firm: 0	Firm Sales: 0	103 local 529	iso 0
Reserves: 0	Non-Firm Sales: 0	- load 50 Copper 69	385
Total Generation for Load: 1135	Total Load Next Hour: 1056	53 598	(If total goes over 600 EPE has to wheel back)
Enter Total Spin Required: 78	PNM Contract: 46	(Contingent upon units 7 & 8 number automatically feeds from the calculation tab)	
Spin Required: 39	IID Firm Contract: 100	Enter Blue Numbers	
Non-Spin Required: 39	TNP Contract: 25	Lambda = \$47.74	
Spin Required + Regulating Margin: 69	SPS Contract: 100	System avg. = \$25.70	
*Amount of Spin: 79	Weighted Avg. Purchase Power Calculator		
Amount of Non-Spin: 50	MWH	\$/MWH	
Total Spin: 129	Firm Block 1: 0	\$0.00	\$0.00
Spin Available/(Deficient): 40	Firm Block 2: 0	\$0.00	\$0.00
Enter Firm Price: \$0.00	Firm Block 3: 0	\$0.00	\$0.00
Enter Non-Firm Price: \$0.00	Firm Block 4: 0	\$0.00	\$0.00
MWH of Firm Avail./(Deficient): 40	Firm Block 5: 0	\$0.00	\$0.00
Total Cost of Firm: \$0	Total: 0	NA	\$0.00 \$0.00
MWH of Non-Firm Avail./(Deficient): 40	MWH	\$/MWH	
Total Cost of Non-Firm: \$0	Non-Firm Block 1: 0	\$0.00	\$0.00
	Non-Firm Block 2: 0	\$0.00	\$0.00
	Non-Firm Block 3: 0	\$0.00	\$0.00
	Non-Firm Block 4: 0	\$0.00	\$0.00
	Non-Firm Block 5: 0	\$0.00	\$0.00
	Total: 0	NA	\$0.00 \$0.00
NOTE: * ACTUAL SPIN SHOW MAY BE LESS SINCE UNIT RAMP RATES ARE NOT CONSIDERED.			\$0.00

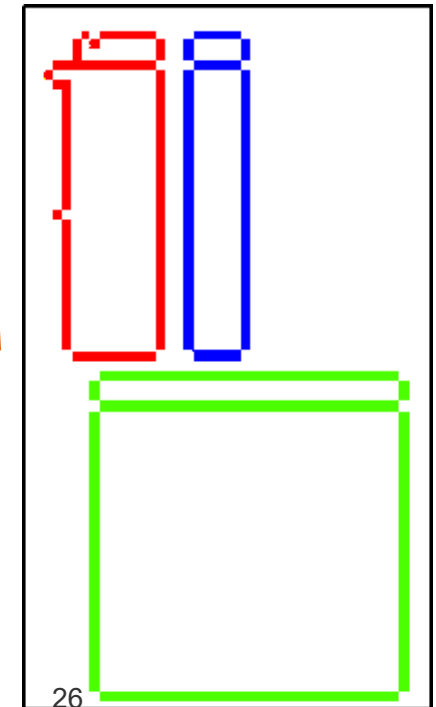
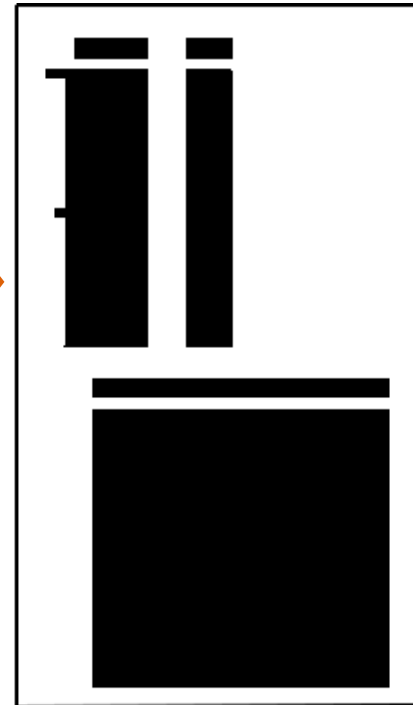
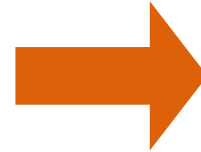
Spin from Ramp rates					
	Output	Highs'	Spin	RR	unloaded
Unit 1	50	80	30	3	30
2	70	82	12	4.5	12
3	0	0	0	5.1	0
GT1	0	0	0	10	0
GT2	0	0	0	10	0
GT1S	0	0	0	3.33	0
GT2S	0	0	0	3.6	0
NM4	148	214	36	3.6	66
Copper	0	0	0	10	0
6	0	0	0	2	0
7	33	33	0	2.1	0
8	96	120	10	1	24
Total	397	529	88		132
FC	49	108			
PV	532	581		Lost Gen.	44
	581	689			

Felix Naumann
Data Quality

Mondrian: Clustering-based Table Recognition

		Maximum Capacity	Change	MTD Avg feb-02	Month-3 Avg ott-01	Month-4 Avg set-01	Tue 05-feb	Mon 04-feb	Sun 03-feb	Sat 02-feb										
Henry Hub	Receipts HH	ACADIAN	200	0	0	0	0	0	0	0										
		BRIDGELINE	80	0	7	32,972	27,804	7	7	7										
		COLUMBIA GU	100	0	0	5,785	4,934	0	0	0										
		DIGCO	0	0	0	0	0	0	0	0										
		JEFFERSON ISE	250	0	7,5	22,367	31,201	7,5	7,5	7,5										
		GULF SOUTH	400	0	75,634	147,07	82,277	75,634	75,634	75,634										
		MAINLINE	180	0	100,733	117,177	120,093	100,733	100,733	100,733										
		HGR	300	0	105,538	81,239	61,301	105,538	105,538	105,538										
		SONAT	125	0	5	0	0	5	5	5										
		SEA ROBIN	250	0	123,367	89,021	157,263	123,367	123,367	123,367										
		TEXAS GAS	0	0	0	0	0	0	0	0										
		TRUNKLINE	75	0	0	10,672	25,791	0	0	0										
		TRANSCO	0	0	0	0	0	0	0	0										
		Total	0	0	-214,246	-303,754	-335,453	-214,246	-214,246	-214,246										

		Thu 31-gen	Wed 30-gen	Tue 29-gen	Mon 28-gen	Sun 27-gen	Sat 26-gen	Fri 25-gen	Thu 24-gen	Wed 23-gen	Tue 22-gen	Mon 21-gen	Sun 20-gen	Sat 19-gen	Fri 18-gen	Thu 17-gen
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
		714	714	714	2,714	2,714	2,714	714	12,714	7,714	8,714	8,714	8,714	8,714	40	0
		1	0	0	1	1	1	1	1	1	1	1	1	1	1	1
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
		201,935	148,085	148,085	1	1	1	0	7,5	5	7,5	7,5	7,5	7,5	7,5	0
		99,067	30,987	30,987	36,513	36,513	36,513	25,906	80,147	115,883	111,402	111,402	111,402	111,402	63,415	118,818
		66,444	119,22	119,22	118,188	118,188	118,188	167,627	172,53	147,527	157,156	157,156	157,156	157,156	170,256	146,999
		78,771	39,196	39,196	256,338	256,338	256,338	207,066	200	200	104,436	104,436	104,436	104,436	140,964	177,797
		87	0	0	0	0	0	0	0	899	0	0	0	0	19,97	0
		30,273	100,268	100,268	91,544	91,544	91,544	114,249	89,961	90,431	93,88	93,88	93,88	93,88	78,357	84,397
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
		0	0	0	0	0	0	0	35,6	39,101	11,5	11,5	11,5	11,5	12	34
		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
		437,292	438,47	438,47	506,298	506,298	506,298	515,563	578,453	606,356	494,589	494,589	494,589	494,589	531,463	561,506
		-51,963	-30,018	-30,018	0	0	0	0	0	-9,323	-21,323	-21,323	-21,323	-21,323	-16,323	-17,823



1. Render spreadsheet as image
2. Recognize elements
3. Cluster elements into tables

ExtractTable: Bad Files – Worse Files

```

min    max    num    dist    mean    std    comment
1.8    1.8    1      0      1.5    0      N
20     60     40     1      40     15     cab
0      0      1      0      # OBIA4RTM config file for setting up Prospect4SAIL
0      1      10     2      #
0.01   0.01   1      0      # Typical values (taken from J Gomez-Dans on https://pypi.org/project/prosail/)
0.009  0.009  1      0      #
0.2    7      40     1      #
-0.35  -0.35  1      2      #
-0.15  -0.15  1      0      #
0.5    0.5    1      0      #
0.2    0.2    1      0      #
0.01   0.01   1      0      #
27.947 27.947 1      0      #
7.04345 7.04345 1      0      #
146.691 146.691 1      0      #
1      1      1      0      #
#
# Parameter | Description of parameter | Units | Typical min | Typical max
#-----|-----|-----|-----|-----
# N | Leaf structure parameter | N/A | 0.8 | 2.5
# cab | Chlorophyll a+b concentration | ug/cm2 | 0 | 80
# caw | Equivalent water thickness | cm | 0 | 200
# car | Carotenoid concentration | ug/cm2 | 0 | 20
# cbrown | Brown pigment | NA | 0 | 1
# cm | Dry matter content | g/cm2 | 0 | 200
# lai | Leaf Area Index | N/A | 0 | 10
# lidfa | Leaf angle distribution | N/A | - | -
# lidfb | Leaf angle distribution | N/A | - | -
# psoil | Dry/Wet soil factor | N/A | 0 | 1
# rsoil | Soil brightness factor | N/A | - | -
# hspot | Hotspot parameter | N/A | - | -
# tts | Solar zenith angle | deg | 0 | 90
# tto | Observer zenith angle | deg | 0 | 90
# phi | Relative azimuth angle | deg | 0 | 360
# typelidf | Leaf angle distribution type | Integer | - | -
#
#
# You can enter your values below -> make sure not to alter the overall structure of this
# template -> otherwise bad things might happen
#
# Further Explanations:
#
# min: Minimum Value of Parameter
# max: Maximum Value of Parameter (in case min=max, the parameter will not be retrieved)
# num: in case min!=max, the number of samples to be drawn for the specific parameter

```

Felix Naumann
Data Quality

Strudel: Verbose CSV Files

Arrest Table		Header				
Arrests for Drug Abuse Violations		Header				
Percent Distribution by Region, 2007		Header				
Drug abuse violations		United States total	Northeast	Midwest	South	West
Total1		100	100	100	100	100
Sale/Manufacturing:	Total	17.5	22.5	18.3	17.1	15
	Heroin or cocaine and their derivatives	7.9	14.2	6.2	7.9	5.5
	Marijuana	5.3	5.7	7.7	4.6	4.7
	Synthetic or manufactured drugs	1.5	1.1	1.1	2.6	0.7
	Other dangerous nonnarcotic drugs	2.8	1.6	3.3	2	4.2
Possession:	Total	82.5	77.5	81.7	82.9	85
	Heroin or cocaine and their derivatives	21.5	22.3	14.7	22.8	22.7
	Marijuana	42.1	44.2	53.1	47.9	29.6
	Synthetic or manufactured drugs	3.3	2.3	3.2	4.3	2.8
	Other dangerous nonnarcotic drugs	15.6	8.6	10.7	7.8	29.9

1 Because of rounding, the percentages may not add to 100.

Metadata
Header
Group header
Data
Aggregation
Notes

Felix Naumann
Data Quality

Strudel: Structure Diversity in Verbose CSV Files

Table 1: Change in UK exports by destination in 2007 (in %)

	2006:1	2006:2	2006:3	2006:4
Germany	11.6	2.8	0.9	2.8
France	11.7	-1.5	-1.9	-1.7
Italy	10.2	-1.2	-2.7	-1.3
Spain	10.7	3.3	13.3	3.1
South Korea	10.9	-3.3	-1.7	-1.4
China	7.9	-1.8	12.2	-1.7
Sweden	11.6	7.8	9.2	7.2
Rest of EU27	12.9	11.9	10.4	10
Rest of Europe	3.1	3.1	7.5	8.9
South America	4.3	8.8	6.4	11.1
Other OPEC	4.2	7.9	7.1	7.8
Rest of the world	1.8	6.7	7.7	6.4
UNEP27	15.5	2.9	12.9	-1.1

Table 2: Percent Distribution by Region, 2007

	Total	United States total	Europe	South America	Other
Pharmaceuticals	100	100	100	100	100
Biotech	13.5	22.5	18.3	17.1	15
Medical devices	7.8	14.2	6.2	7.5	5.5
Medical equipment	5.3	5.7	7.7	4.6	4.7
Biotech or pharmaceutical drugs	1.5	1.1	1.1	2.6	8.7
Pharmaceutical drugs	2.8	1.6	3.3	7	4.2
Pharmaceuticals	82.5	77.5	81.7	82.9	80
Medical devices	21.5	22.3	14.7	22.8	22.7
Medical equipment	42.1	46.2	51.1	47.9	25.6
Biotech or pharmaceutical drugs	3.3	2.3	3.2	4.3	2.8
Pharmaceutical drugs	15.4	8.6	18.7	7.8	25.9

Table 3: Detailed breakdown of pharmaceuticals

Region	Pharmaceuticals	Medical devices	Medical equipment	Biotech or pharmaceutical drugs	Pharmaceutical drugs
United States	100	100	100	100	100
Europe	100	100	100	100	100
South America	100	100	100	100	100
Other	100	100	100	100	100

Table 4: Summary of data aggregation

Category	Value	Sub-category	Value	Sub-category	Value
Pharmaceuticals	100	Medical devices	10	Medical equipment	10
Biotech or pharmaceutical drugs	10	Pharmaceutical drugs	10		

Table 5: Detailed breakdown of medical equipment

Region	Medical equipment	Biotech or pharmaceutical drugs	Pharmaceutical drugs
United States	100	100	100
Europe	100	100	100
South America	100	100	100
Other	100	100	100

Table 6: Summary of data aggregation

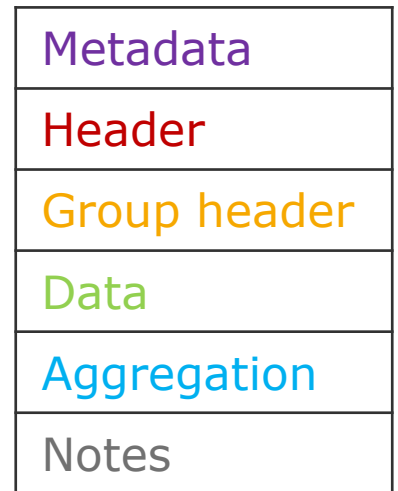
Category	Value	Sub-category	Value	Sub-category	Value
Medical equipment	100	Biotech or pharmaceutical drugs	10	Pharmaceutical drugs	10

Table 7: Detailed breakdown of pharmaceuticals

Region	Pharmaceuticals	Medical devices	Medical equipment	Biotech or pharmaceutical drugs	Pharmaceutical drugs
United States	100	100	100	100	100
Europe	100	100	100	100	100
South America	100	100	100	100	100
Other	100	100	100	100	100

Table 8: Summary of data aggregation

Category	Value	Sub-category	Value	Sub-category	Value
Pharmaceuticals	100	Medical devices	10	Medical equipment	10
Biotech or pharmaceutical drugs	10	Pharmaceutical drugs	10		



Felix Naumann
Data Quality

AggreCol: Aggregations in CSV Files

		$\% \text{ Change 2003 vs. 2002} = \frac{\text{FY2003} - \text{FY2002}}{\text{FY2002}}$							
Income Statement Data	% Change	2003 vs. 2002	FY2003	FY2002	FY2001	FY2000	FY1999	FY1998	FY1997
Hardware Revenue	2.2%	\$137,013	\$134,121	\$116,058	\$152,186	\$155,237	\$126,974	\$102,816	
Software Revenue	17.8%	\$71,251	\$60,484	\$55,873	\$66,290	\$63,317	\$57,744	\$45,985	
Service Revenue	11.2%	\$191,927	\$172,558	\$154,845	\$143,378	\$118,525	\$97,200	\$79,368	
Total Revenue	9.0%	\$400,191	\$367,163	\$326,776	\$361,854	\$337,079	\$281,918	\$228,169	
Memo Item:									
Maintenance Revenue (included in Service Revenue)	13.9%	\$113,274	\$99,467	\$87,007	\$65,628	\$54,953	\$45,908	\$37,38	
Hardware Gross Profit	2.3%	\$38,977	\$38,116	\$40,683	\$51,462	\$50,670	\$43,947	\$39,267	
Hardware Gross Profit %	-	28.4%	28.4%	35.1%	33.8%	32.6%	34.6%	38.2%	
Software Gross Profit	11.5%	\$54,045	\$48,457	\$46,875	\$51,349	\$52,138	\$47,235	\$37,464	
Software Gross Profit %	-4.2 Points	75.9%	80.1%	83.9%	77.5%	82.3%	81.8%	81.5%	
Service Gross Profit	16.5%	\$105,538	\$90,564	\$76,472	\$71,741	\$61,367	\$46,455	\$39,447	
Service Gross Profit %	+2.5 Points	55.0%	52.5%	49.4%	50.0%	51.8%	47.8%	49.7%	
Total Gross Profit	12.1%	\$198,560	\$177,137	\$164,028	\$174,552	\$164,175	\$137,637	\$116,178	
Gross Profit %	+1.4 Points	49.6%	48.2%	50.2%	48.2%		48.8%	50.9%	

$$\text{Hardware GP \%} = \frac{\text{Hardware GP}}{\text{Hardware Revenue}}$$

$$\text{Service GP \%} = \text{FY2003} - \text{FY2002}$$

$$\text{Total Gross Profit} = \text{Hardware GP} + \text{Software GP} + \text{Service GP}$$

Felix Naumann
Data Quality

Ill-formed Records Abort Data Loading

1	NAME	LAST NAME	JOBTITLE	DESCR	HIRE_DT	ANNUAL_RT	GROSS
2	Aaron	Kareem D	Utilities Inst Repair I	A50550	08/27/2018	32470	25743.94
3	Aaron	Patricia G	Office Services II	A03031	10/24/1979	60200	57806.13
4	Abadir	Adam O	Council Technician	A02002	12/12/2016	64823	64774.11
5	Abaku	Aigbolosimuan O	Police Officer	A99094	04/17/2018	53640	59361.55
6	Abbeduto	Mack	Assistant State Attorney	A29011	05/22/2017	68562	61693.59
7	Abbott	Ethan N	Recreation Arts Instructor	A68002	04/11/2018	33280	26156.48
8	AbbottCole	Michelle	Operations Officer III	A90005	11/28/2014	75110	75529.99
9	Abdal Rahim	Naim A	Fire Pump Operator Suppression	A64120	03/30/2011	69595	82132.61
10	Abdi	Ezekiel W	Police Sergeant	A99160	06/14/2007	93284	122992.1
11	Abdul Adl	Attrice A	Radio Dispatcher Sheriff	A38410	09/02/1999	50079	58459.78
12	Abdul Aziz	Hajr E	Swimming Pool Operator	P04002	06/01/2017	28554	15807.64
13	Abdul Aziz	Yaqub M	Swimming Pool Operator	P04002	06/01/2017	28554	6417.9
14	Abdul Saboor	Dana N	Paralegal	A99393	04/13/1998	57857	35715.2
15	Abdul	Jalil	Engineer I	A50101	07/17/2017	64505	65577.5
16	Abdul-Jabbar	Bushra A	Social Service Coordinator	A65028	04/14/2008	46395	46359.55
17	Abdul-Khaliq	Amahl	Recreation Leader II	A04005	06/06/2019	32131	5344.5
18	Abdullah	Beverly A	Office Support Specialist III	A06004	12/01/1986	41757	44522.96
19	Abdullahi	Sharon M	911 Operator	A64604	Wednesday, 6 October 2004	56322	53751.88
20	Abdullateef	Muhammed L	Supt of Public Bldg Repair	A85001	05/09/2019	78000	9000
21	Abdulrahman	Mustafa H	Police Officer Trainee	A99416	12/28/2018	53512	25347.99
22	Abdul Saboor	Jamillah	Printer Library	A75055	07/27/2009	44584	42555.83
23	Abdunafi	Karim	Community Aide	A04015	06/13/2019	24960	9888.55
24	Abdur-Rahman	Diane	Office Services I	A03092	03/27/2017	25363	25649.94
25	Abdurrahman	Saleh Z	Lifeguard I	P04002	06/04/2019	23920	1245.5
26	Abebe	Miraf E	Auditor II	A24002	Saturday, 2 June 2012	67236	60600.92
27	Abend Kollin	Emily L	Fleet Quality Control Analyst	A85301	01/05/2017	47828	49133.1133
28	Abid	Amal	Engineer II	A49102	12/02/2013	71774	71630.83
29	Abid	Paula	Recreation Arts Instructor	A04009	06/25/2007	24960	3591.96
30	Abid	Paula	Recreation Arts Instructor	A04009	06/25/2007	24960	3591.96

Suragh: Row Patterns – Outlier Rows

SFY, Fund_Source, Age, Total Children Receiving CHDP Services, Total, ...
 2008-2009, All, 0, 557757, 24.34%, "\$21,840,767 ", 38.03%
 2008-2009, All, 1, 314994, 13.75%, "\$7,262,306 ", 12.64%

...

2008-2009, All, 10, 55674, 2.43%, "\$1,146,066 ", 2.00% SFY, Fund_Source, ...
 2008-2009, All, 0, 557757, 24.34%, "\$21,840,767 ", 38.03%
 2008-2009, All, 1, 314994, 13.75%, "\$7,262,306 ", 12.64%

...

2008-2009, All, Unknown, 1063, 0.05%, \$0 , 0.00%
 2008-2009, All, TOTAL, 2291689, 100.00%, "\$57,436,517 ", 100.00%
 2008-2009, FFS, 0, 169699, 41.75%, "\$14,450,511 ", 45.16%
 2008-2009, FFS, 1, 48823, 12.01%, "\$5,157,938 ", 16.12%
 2008-2009, FFS, 2, 29241, 7.19%, "\$2,019,526 ", 6.31%

...

2008-2009, FFS, UNKNOWN, 0, 0.00%, \$0 , 0.00%
 2008-2009, FFS, TOTAL, 406504, 100.00%, "\$31,995,243 ", 100.00%
 2008-2009, GFS, 0, 91612, 33.09%, "\$7,363,856 ", 31.33%
 2008-2009, GFS, 1, 22723, 8.21%, "\$2,061,059 ", 8.77%

...

2008-2009	All	<SEQD>	<SEQD>	<SEQD>.<SEQD>%	"\$<NUM> "	<SEQD>.<SEQD>%
2008-2009		<SEQD>	<SEQD>	<SEQD>.<SEQD>%	"\$<NUM> "	<SEQD>.<SEQD>%

Felix Naumann
Data Quality

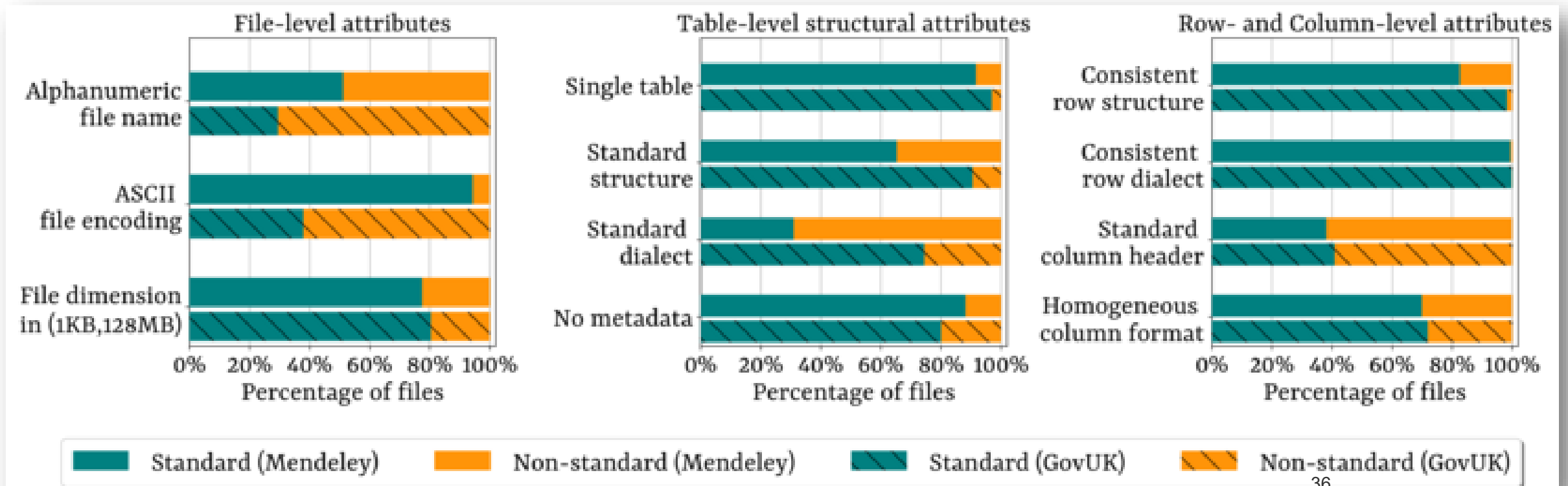
Pollock: Benchmarking the Ingestion Ability of Systems

```
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
```

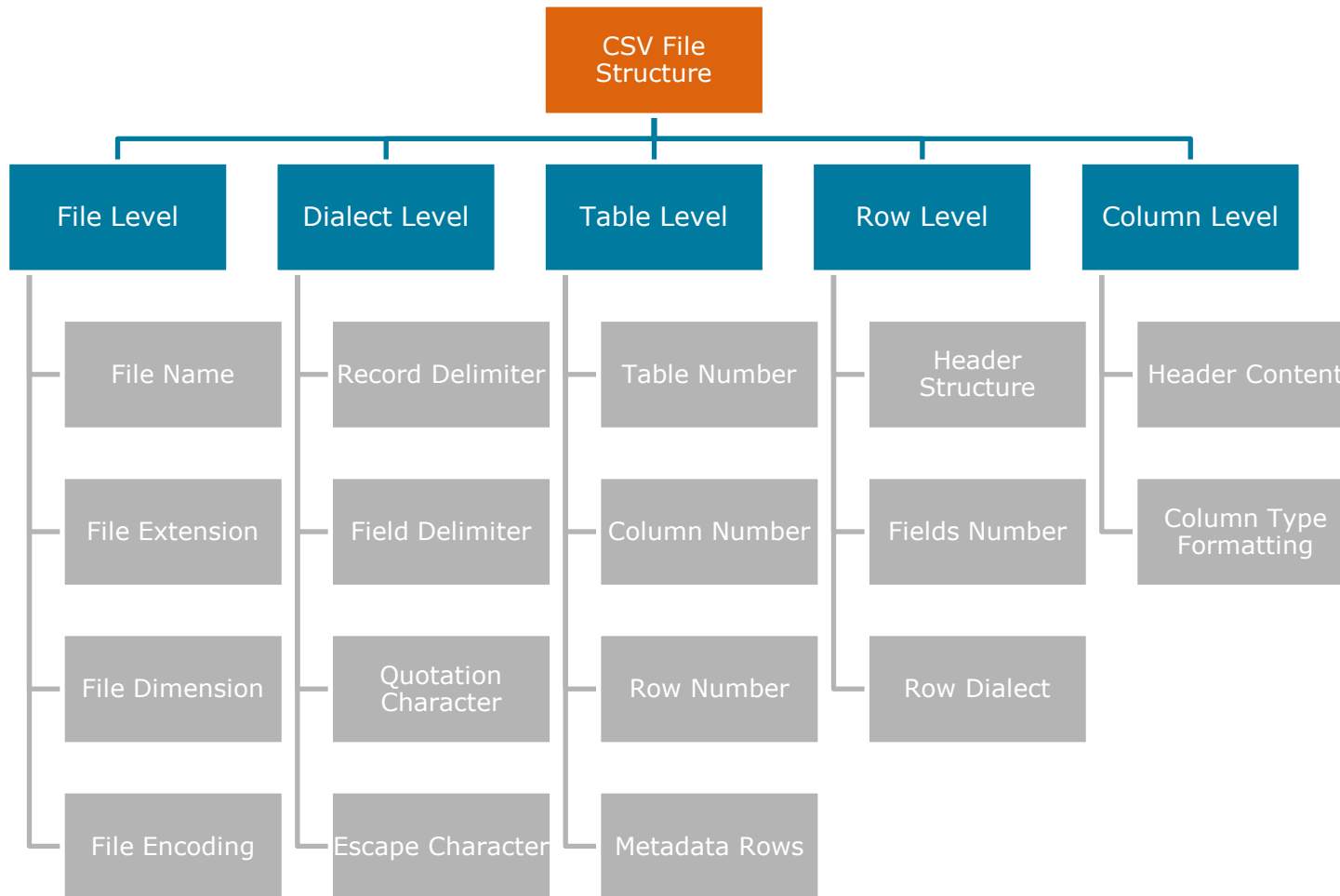
Felix Naumann
Data Quality

Pollock: Raw Data Survey

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



Pollock: Benchmark Dimensions and Results



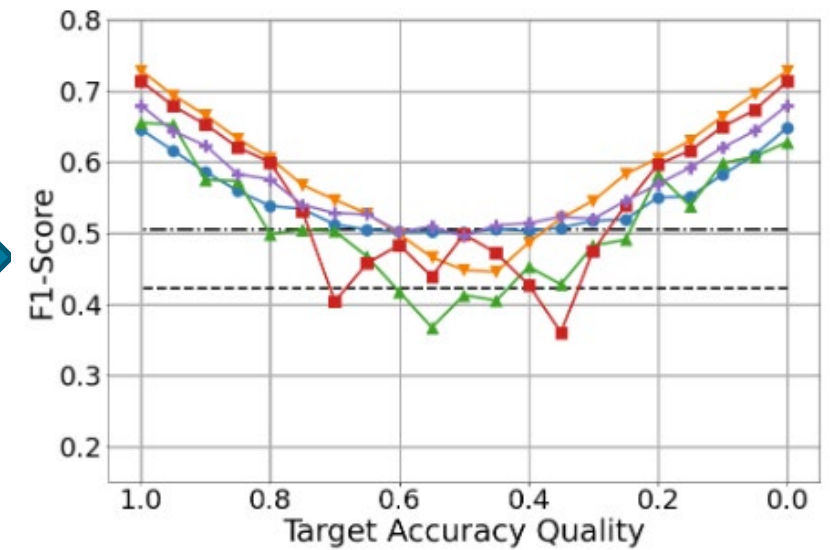
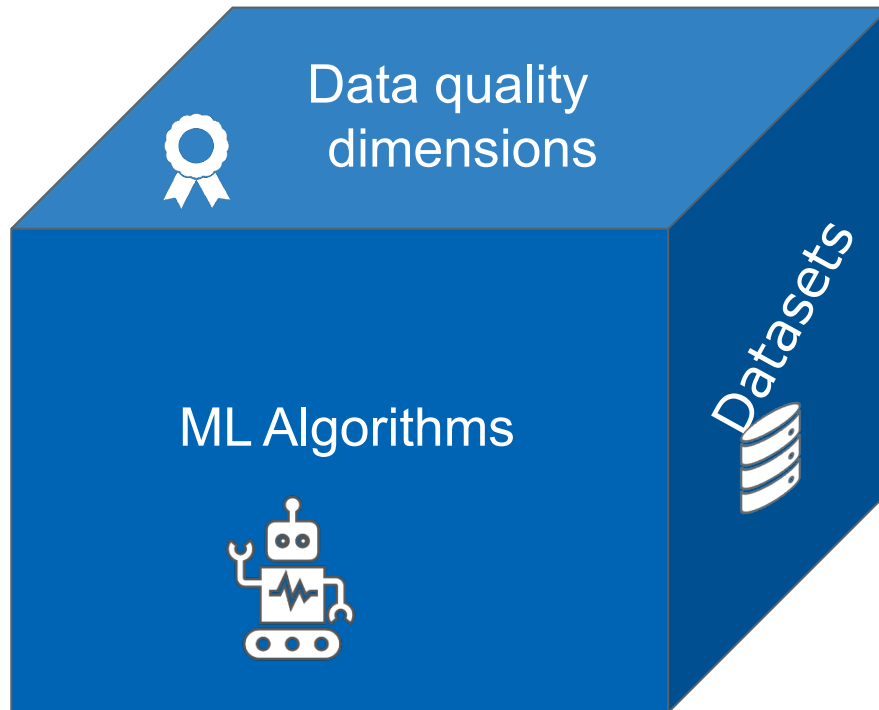
	Pollock score (2 289 +1) files	
	Simple	Weighted
CLEVERCSV 0.7.4	9.05	9.49
CSVCOMMONS 1.9.0	6.63	9.29
HYPOPARSR 0.1.0	3.73	4.41
OPENCsv 5.6	6.62	7.80
PANDAS 1.4.3	9.88	9.75
PyCsv 3.10.5	9.71	9.47
RCsv 4.2.1	7.78	6.76
UNIVOCITY 2.9.1	9.35	7.97
MARIADB 10.9.3	8.81	7.44
MySQL 8.0.31	8.88	7.45
POSTGRESQL 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

Runs

- 5 runs, average

Datasets

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

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

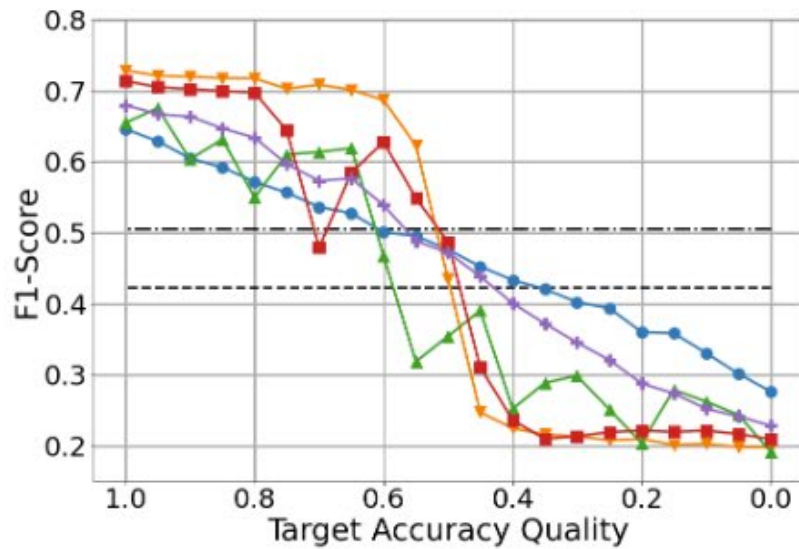
Scenarios

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

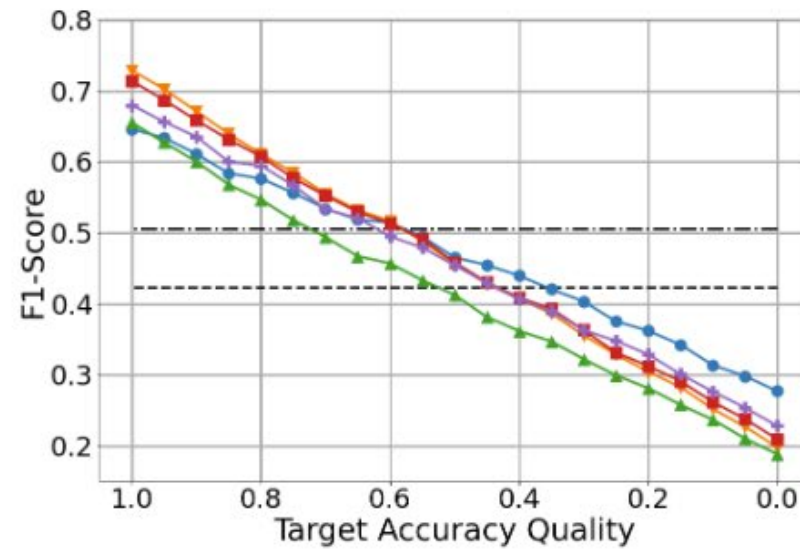
Felix Naumann
Data Quality

Together: 3,780 experiments

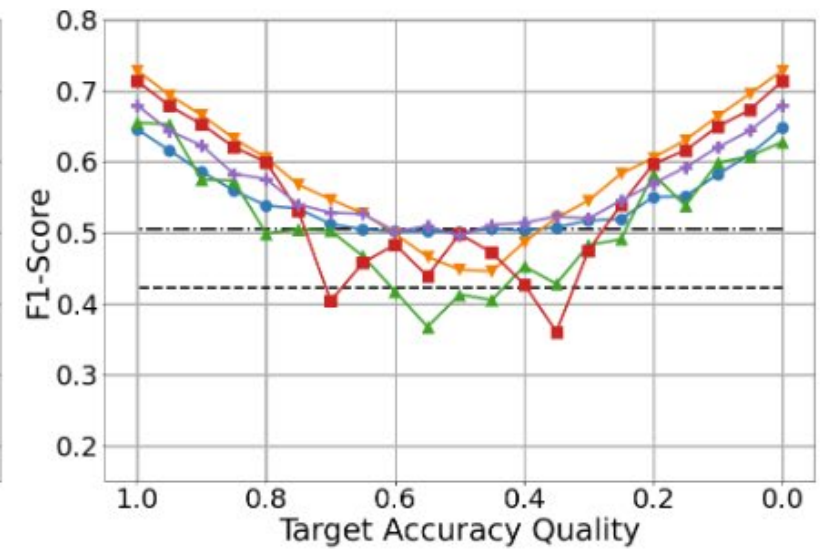
Example Results



Pollute Training



Pollute Test



Pollute Training & Test

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

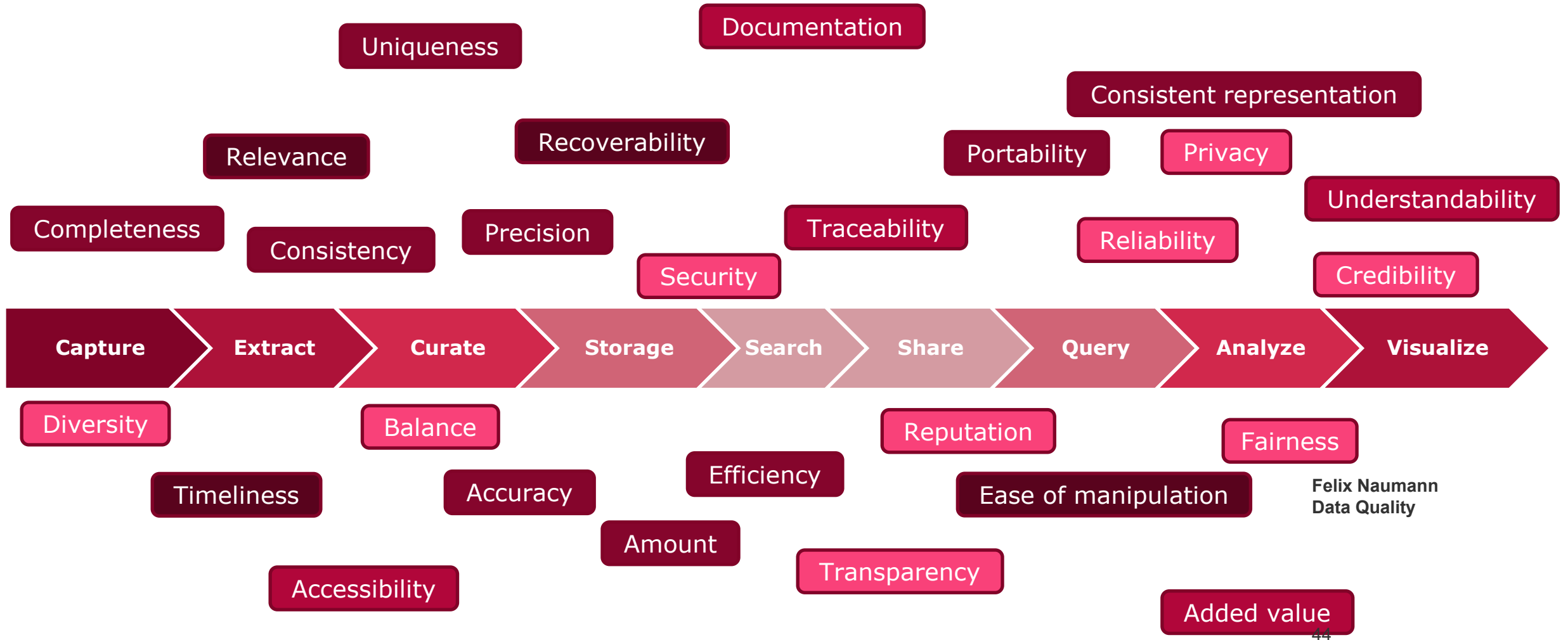
- Majority class baseline
- Class ratio baseline
- Decision Tree
- Logistic Regression
- Multilayer Perceptron
- Support Vector Machine
- k-Nearest Neighbors

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

Training, validation and testing data sets shall be **relevant**, sufficiently **representative**, and to the best extent possible, **free of errors** and **complete** in view of the intended **purpose**. They shall have the **appropriate statistical properties**, including, where applicable, as regards the **persons** or groups of persons in relation to whom the high-risk AI system is intended to be used.



Data Quality along the AI Pipeline



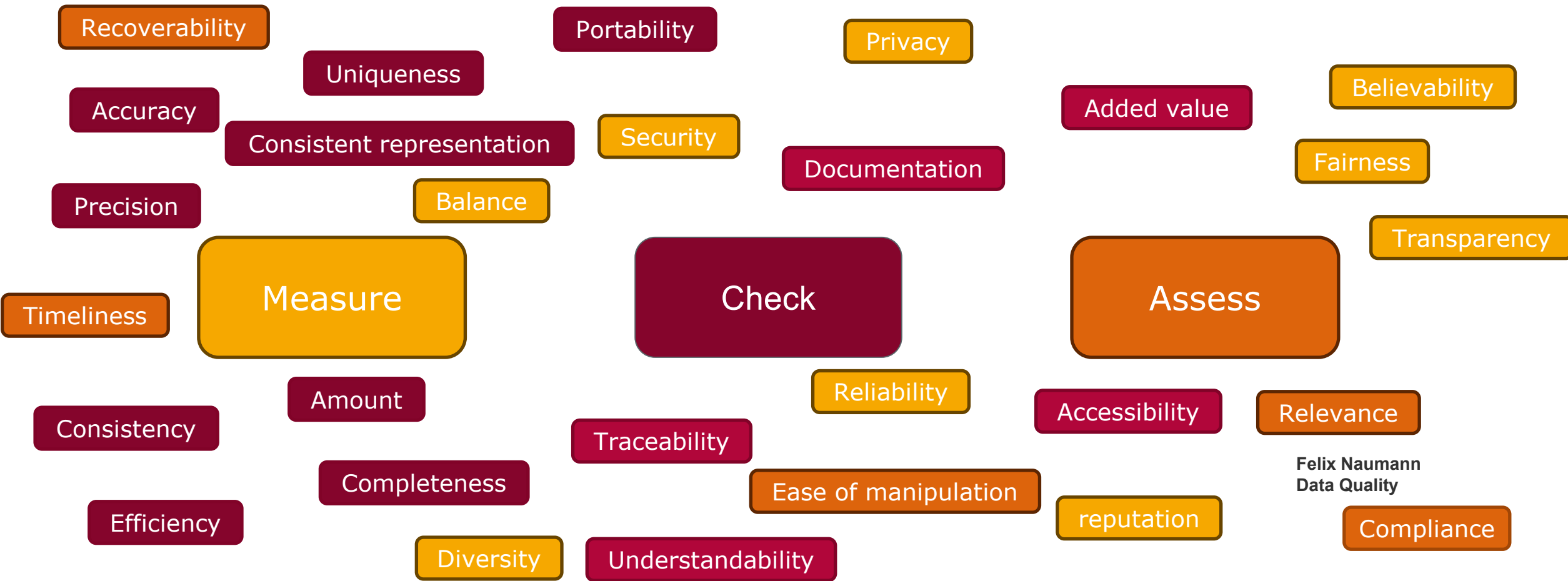
Felix Naumann
Data Quality

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

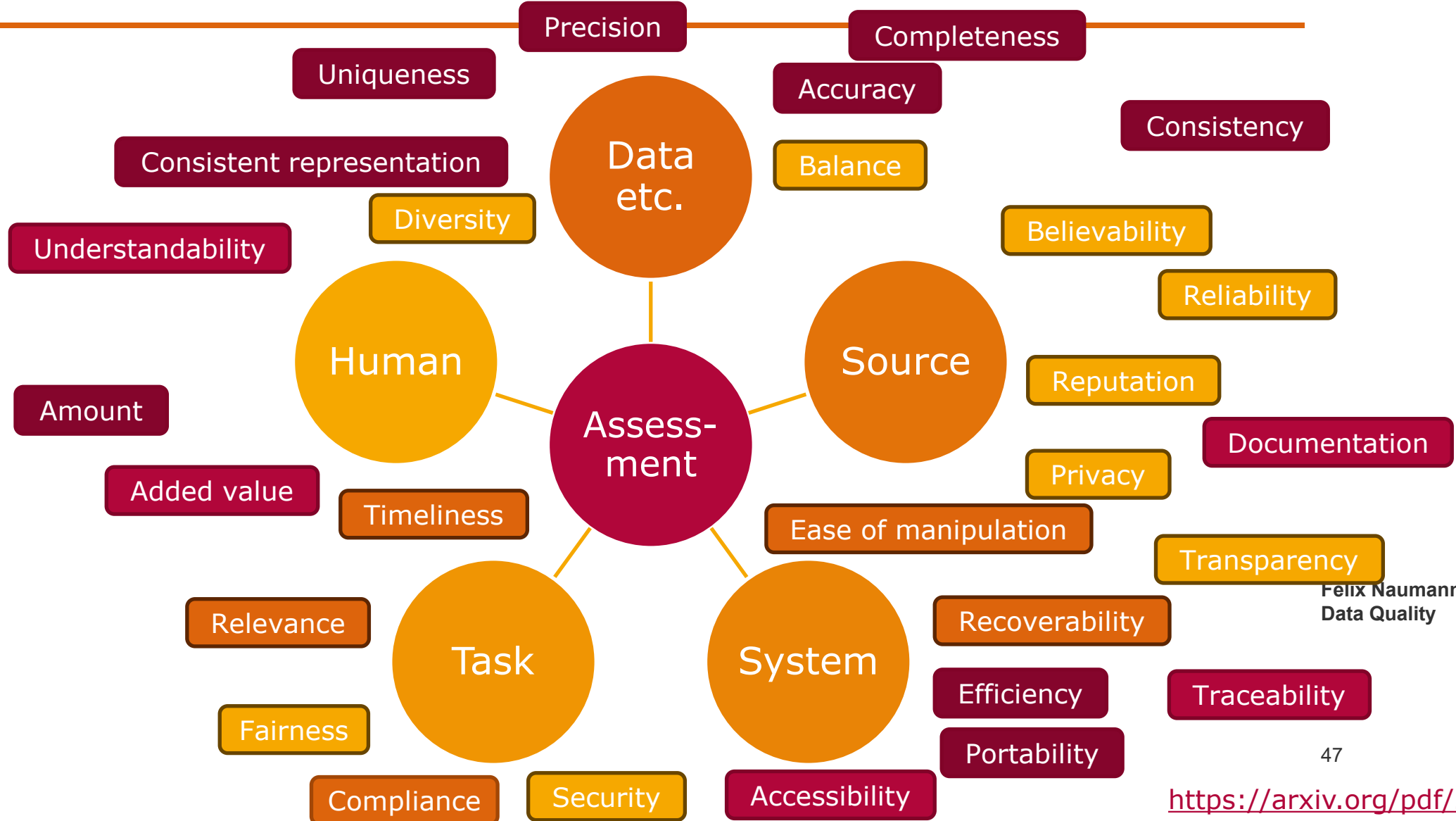


Assessing Data Quality



Felix Naumann
Data Quality

Ingredients for DQ Assessment: Five Facets



Felix Naumann
Data Quality

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

