

# Von Apache Lucene zu Multi-Model Data Stores (Document, Key-Value, Graph) mit multidimensionaler bitemporaler 6NF Ablage + Dynamische Links



8. Tagung (Herbsttagung 2018)





**42**  
data **mor**row



**data**  
**m**odeling  
**zone**

**EUROPE2018**

24<sup>th</sup> & 25<sup>th</sup> September 2018  
**Dusseldorf, Germany**



# Agenda



## • Fuzzy / Roter Faden?

DISCLAIMER:

Denn das ist alles nur geklaut

Das ist alles gar nicht meine

Das ist alles nur geklaut

Doch das weiß ich nur ganz alleine

Das ist alles nur geklaut

Und gestohlen

Nur gezogen

Und geraubt

Entschuldigung, das hab' ich mir erlaubt

Aus „Alles nur geklaut“: Die Prinzen

# KI ist eindeutig das neue Schlagwort - die Data Fashion Show geht weiter.



- Letztes Jahr verschwand das Schlagwort Big Data endgültig, aber alle sprachen von Data Science. In diesem Jahr sprachen alle über KI und maschinelles Lernen. Ich schätze, es ist ein Teil der Zugehörigkeit zum Datenberuf, sich jedes Jahr an das neueste Schlagwort anzupassen! Meine Keynote wurde vom Redaktionsleiter der Konferenz von "data science for grown ups" in "succeeding in AI" umbenannt. In meinem eigenen Unternehmen war es schwierig, den Leuten zu erklären, dass unsere Informatiker auch an KI arbeiten können (im Grunde genommen habe ich in den meisten Präsentationen KI geschrieben). Ich habe mich von nun an entschieden, Datenwissenschaftler die Daten- und KI-Wissenschaftler zu nennen, um zu signalisieren, dass es sich um die gleichen Personen handelt. Aus kritischer Sicht könnte man argumentieren, dass die echte KI noch etwas Zeit in Anspruch nimmt. Die Art von KI, über die die meisten Leute sprechen, ist nicht ein einzelner Algorithmus, sondern ein System einer Vielzahl von miteinander verbundenen Algorithmen, die in eine gut verwaltete Algorithmusarchitektur eingebettet sind, die Maschinen sinnvoll machen kann, denken, lernen und handeln (siehe auch meinen letzten Artikel über Smart Machines). Soweit ich weiß, hat das bisher noch niemand erreicht. Mein Kollege Dat Tran, Head of Data Science bei Idealo, hat es heute in einem LinkedIn-Posting gut zusammengefasst: "Wir sind so weit entfernt von dem, was die Leute echte KI nennen. Alles, was wir tun, ist eine mathematische Optimierung einiger Daten. Und wenn die Daten schlecht sind, ist auch die KI schlecht". Eine Sache, auf die ich in meiner heutigen Keynote hingewiesen habe, war, dass wir es versäumt haben, eine Wissenstaxonomie für unsere Maschinen zu erstellen, sie alle müssen sich auf Rohdaten verlassen, wenn wir maschinelle Lernalgorithmen ausführen. Aber wie können Maschinen etwas über die Gesellschaft lernen, wenn sie keine Bücher lesen?
- Übersetzt mit [www.DeepL.com/Translator](http://www.DeepL.com/Translator)

Dr. Alexander Borek Global Head of Data & Analytics at Volkswagen Financial Services AG  
Succeeding in Data & AI: Insights from the Data Leaders Summit Europe 2018  
Veröffentlicht am 19. Oktober 2018 auf LinkedIn

# Der Anfang: Herodot 1,26-33 Kroisos und Solon



„Wenn du den Halys überschreitest, wirst du ein großes Reich zerstören.“

(1,26) Τελευτήσαντος δὲ Ἀλυάττεω ἐξεδέξατο τὴν βασιληίην Κροῖσος ὁ Ἀλυάττεω, ἐτέων ἑὼν ἡλικίην πέντε καὶ τριήκοντα, ὃς δὴ Ἑλλήνων πρῶτοισι ἐπεθήκατο Ἐφεσίοισι. Ἐνθα δὴ οἱ Ἐφέσιοι πολιορκεόμενοι ὑπ' αὐτοῦ ἀνέθεσαν τὴν πόλιν τῇ Ἀρτέμιδι, ἐξάψαντες ἐκ τοῦ νηοῦ σχοινίον ἐς τὸ τεῖχος· ἔστι δὲ μεταξὺ τῆς τε παλαιῆς πόλιος, ἣ τότε ἐπολιορκέετο, καὶ τοῦ νηοῦ ἑπτὰ στάδιοι. Πρῶτοισι μὲν δὴ τούτοισι ἐπεχείρησε ὁ Κροῖσος, μετὰ δὲ ἐν μέρει ἑκάστοισι Ἴώνων τε καὶ Αἰολέων, ἄλλοισι ἄλλας αἰτίας ἐπιφέρων, τῶν μὲν ἐδύνατο μέζονας παρευρίσκειν, μέζονα ἐπαιτιώμενος, τοῖσι δὲ αὐτῶν καὶ φλαῦρα ἐπιφέρων.

Nach dem Tod des **Alyattes** übernahm dessen Sohn **Kroisos** die Herrschaft, der bereits in einem Alter von 35 Jahren stand und unter allen **Hellenen** zuerst die **Ephesier** angriff. Als nun die **Ephesier** von ihm belagert wurden, weihten sie ihre Stadt der **Artemis**, nachdem sie von dem Tempel bis zur Stadtmauer ein Seil gebunden hatten; es liegen aber zwischen der alten Stadt, die damals belagert wurde, und dem Tempel sieben Stadien. Diese also griff **Kroisos** zuerst an, nachher aber griff er die übrigen **Ionier** und **Aioler** der Reihe nach an, indem er bei jeder Stadt einen anderen Grund vorgab, einen erheblicheren, wenn er einen solchen ausfindig machen konnte; bei einigen von ihnen nahm er auch einen ganz geringen Vorwand.

# Analyse der deutschen Übersetzung



- Wörter 118
- Eindeutige Wörter 86
- „die“ ist das am häufigsten vorkommende Wort. 5 x
- Wörter mit Länge größer als 4 = 58
- Wichtige Wörter = ca. 8 – 12
- U.a. Aioler, Alyattes, Artemis, belagert, Ephesier, Hellenen, Herrschaft, Ionier, Kroisos, Stadt, Stadtmauer, Tempel ...

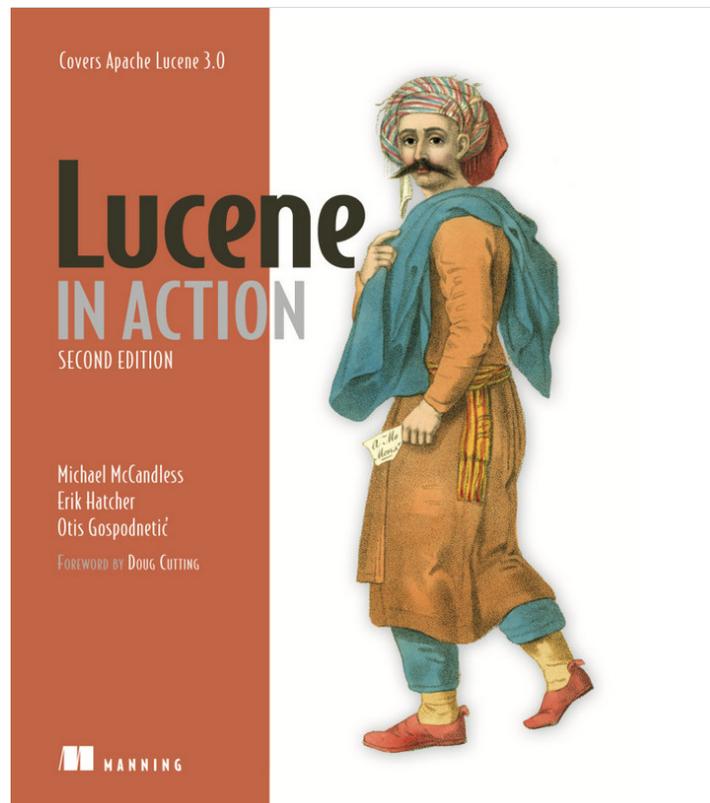
# Lucene



- Vor MongoDB, vor Cassandra, vor "NoSQL" gab es Lucene.
- Wusstest du, dass Doug Cutting 1999 die ersten Versionen von Lucene geschrieben hat?
- Um die Dinge in Zusammenhang zu bringen, war dies etwa zu der Zeit, als Google mehr ein Forschungsprojekt als eine tatsächlich vertrauenswürdige Anwendung war. Googles Proof-of-Concept-Suchmaschine war immer noch ein weitläufiger Satz von Desktop-Computern in Stanfords Forschungslabors.
- Übersetzt mit [www.DeepL.com/Translator](http://www.DeepL.com/Translator)

Lucene: The Good Parts auf Parse.ly by Andrew Montalenti March 12, 2015  
<https://blog.parse.ly/post/1691/lucene/> heruntergeladen am 24.10.2018

# Lucene in Action



- ... Lucene nähert sich den Problemen der Datenexploration aus der Sicht der "Information Retrieval", nicht aus der Sicht der "Database Management Theory". ...
- Lucene's Schöpfer erwägt: Wie unterstützen wir Abfragen, die normale Benutzer tatsächlich eingeben werden? Wie können wir schnell alle Daten, die wir haben, auf einen Schlag durchsuchen? Wie ordnen wir die Ergebnisse an, wenn es mehr als eine wahrscheinliche Übereinstimmung gibt? Wie fassen wir die gesamte Ergebnismenge zusammen, auch wenn wir nur genügend Platz haben, um einen Teil der Ergebnismenge anzuzeigen?
- Damals existierten Solr und Elasticsearch noch nicht.
- Übersetzt mit [www.DeepL.com/Translator](http://www.DeepL.com/Translator)

Lucene: The Good Parts auf Parse.ly by Andrew Montalenti March 12, 2015  
<https://blog.parse.ly/post/1691/lucene/> heruntergeladen am 24.10.2018

Lucene in Action: <https://www.manning.com/books/lucene-in-action-second-edition>

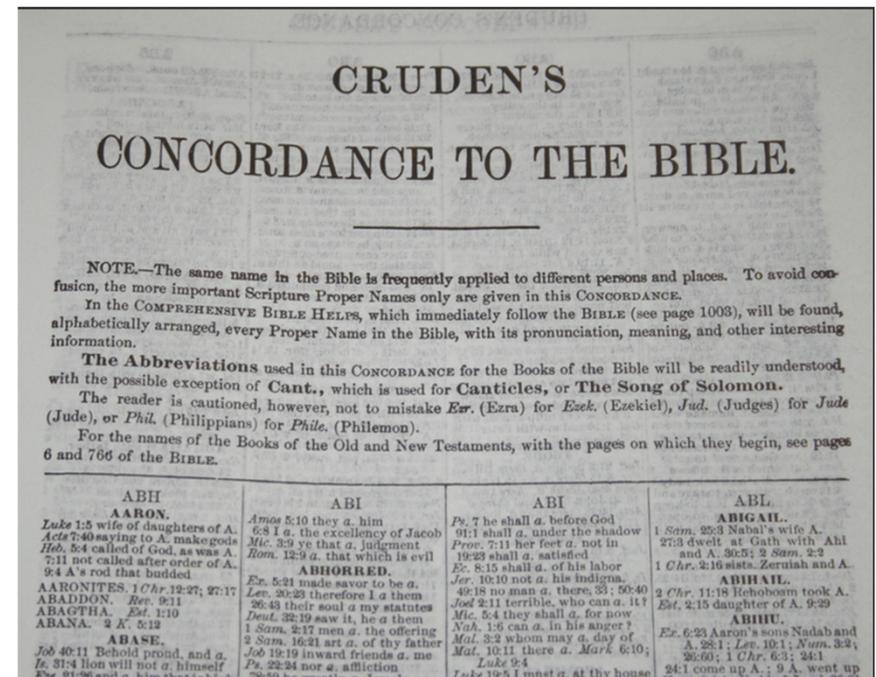
# Vokabular und Corpus



- Vokabular: Der gesamte Satz eindeutiger Begriffe in einem Corpus.
- Corpus: der gesamte Satz von Dokumenten in einem Index

Übrigens, dies ist eine ziemlich alte Technik der Datengewinnung. Das erste vollständige Vokabular eines komplexen Textes wurde im Jahr 1262 von 500 sehr geduldigen Mönchen erstellt. Das fragliche Dokument war natürlich die Bibel, und der Wortschatz wurde als Konkordanz bezeichnet. Wie lautet das Sprichwort? "Es gibt nichts Neues unter der Sonne."

Lucene: The Good Parts auf Parse.ly by Andrew Montalenti March 12, 2015  
<https://blog.parse.ly/post/1691/lucene/> heruntergeladen am 24.10.2018



Ein Beispiel für eine biblische Konkordanz, die etwas moderner ist als die, an der die Mönche gearbeitet haben.

# Begriffserklärungen



- Hier sind einige Begriffe, die Sie in den Lucene und Information Retrieval Communities finden, die in den SQL und Datenbank Communities nicht annähernd so verbreitet sind. Lucene definiert sogar den Begriff " **term** " neu.
- **document**: ein Datensatz; die Sucheinheit; die zurückgegebene Sache als Suchergebnis ("keine Zeile")
- **field**: ein typisierter Slot in einem Dokument zum Speichern und Indizieren von Werten ("keine Spalte").
- **index** : eine Sammlung von Dokumenten, typischerweise mit dem gleichen Schema ("keine Tabelle").
- **corpus**: der gesamte Satz von Dokumenten in einem Index
- **inverted index**: interne Datenstruktur, die Begriffe nach ID auf Dokumente abbildet.
- **term**: aus dem Quelldokument extrahierter Wert, der für den Aufbau des invertierten Index verwendet wird.
- **vocabulary**: der gesamte Satz eindeutiger Begriffe in einem Korpus.
- **uninverted index**: alias "Felddaten", Anordnung aller Feldwerte pro Feld, in Dokumentenreihenfolge
- **doc values**: alternative Art der Speicherung des uninvertierten Index auf der Festplatte (Lucene-spezifisch)

Lucene: The Good Parts auf Parse.ly by Andrew Montalenti March 12, 2015

<https://blog.parse.ly/post/1691/lucene/> heruntergeladen am 24.10.2018

# Umkehrung unseres Corpus



- Beginnen wir mit einem einfachen Corpus von: Zwei Dokumente, doc1 und doc2, enthalten beide das Feld "tag", Typ "string", mit dem Text "big data". Es gibt auch doc3, gleiche Struktur, aber sein Tag enthält den Text „small data“.
- Wie können wir mit diesem kleinen Corpus Dinge finden?
- Anstatt so zu speichern:

```
doc1={"tag": "big data"}  
doc2={"tag": "big data"}  
doc3={"tag": "small data"}
```

# Umkehrung unseres Corpus



- Wir können den "invertierten Index" speichern. Ja, was ist das denn?

```
big=[doc1,doc2]
data=[doc1,doc2,doc3]
small=[doc3]
```

- Ah, also ist es kein Index von Dokumenten zu Begriffen, es ist ein Index von Begriffen zu Dokumenten. Schlaue. Wenn wir die Daten so organisieren, können wir Dokumente schneller nach Wert finden. Wenn ich nach „big“ suche, bekomme ich doc1 und doc2 zurück. Wenn ich nach „small“ suche, bekomme ich doc3 zurück. Wenn ich nach „data“ suche, bekomme ich alle Dokumente zurück. Dies ist im Grunde die Kerndatenstruktur in Lucene und in der Suche im Allgemeinen. Ein Hurra für den invertierten Index!

# Nicht in meinem Vokabular



- In den obigen Dokumenten habe ich 3 "Begriffe", und die Annahme ist, dass ich sie durch grundlegende Whitespace-Tokenisierung generiert habe. So hatte mein ursprünglicher Korpus die Feldwerte [„big data“, „small data“], aber meine generierten Begriffe sind [„big“, „small“, „data“].
- Das deutet bereits auf etwas Interessantes an Begriffen hin. Wenn Informationen in Ihren Feldwerten wiederholt werden, werden sie durch Herausziehen der Begriffe komprimiert.

# Nicht in meinem Vokabular



- Begriffe sind interessant, wenn Sie Daten haben, die sich häufig unter Ihren Dokumenten wiederholen. In diesem kleinen Beispiel wird der Begriff „data“ in beiden Dokumenten wiederholt, erfordert aber nur einen Eintrag im invertierten Index. Stellen Sie sich die gleiche Art von Corpus vor wie oben, aber wenn Sie 1.000 Gesamtdokumente haben, dann ist die Hälfte mit „big data“ und die andere mit „small data“ versehen. In diesem Fall haben Sie vielleicht:

```
data=[1,2,3,...,1000]
big=[1,3,5,7,9,...,999]
small=[2,4,6,8,...,1000]
```

Aus „BIG DATA“ wird „SMALL DATA“

# Buchindex

## Bin ich Krösus, oder was?



<p>328</p> <p style="text-align: center;">Index</p> <p>Canvas of life turned upside down, 68</p> <p>"Carbonate of pork," 325</p> <p>Carracci, the, 147</p> <p><i>Casina di Banda</i>, 119</p> <p>Castelletto, 265</p> <p>Cavagnago, 76</p> <p>Cenere, Monte, narcissuses on, 228</p> <p>Ceres, 161</p> <p><i>Crisca</i>, 133</p> <p>Chalk, Conté, the Italian for whom this was the one thing needful, 136</p> <p>Chalk eggs, 43</p> <p>Chamois, foot of, 283</p> <p>Change, repudiation of desire for sudden, 186</p> <p>— importance of, depends on the rate of introduction, 196</p> <p>— either the circumstances or the sufferer will, 196</p> <p>Changes, sweeping, to be felt hereafter as vibrations, 60</p> <p><i>Chiapissino</i>, 163</p> <p><i>Chiese and the abbi</i>, 289</p> <p>Cherries, 33, 34, 46</p> <p>Chestnuts, 118</p> <p>Chicory and seed onions, weary utterness in, 227</p> <p>Children, subalpine, 301</p> <p>— what becomes of the clever, 149</p> <p>Chinese, the examination-ridden, 151</p> <p>Chronic, 75</p> <p>"Chow," 52</p> <p>Church-going, subalpine, 303</p> <p>Circulation of people like blood, 20</p> <p>Ciseri, his picture at Locarno, 271</p> <p>Civilisation, antiquity of Italian, 124</p> <p>— stationary, of ants and bees, 195</p> <p>Class distinction inevitable, 195</p> <p>Classification only possible through sense of shock, 63</p> <p>Close our English and S. Michele</p>	<p style="text-align: center;">Index</p> <p>329</p> <p>Cooling, Wednesbury, 55, 305</p> <p>Collects, unsympathetic priest bristling with, 111</p> <p>Colleone, Medea, 231</p> <p>Colma di San Giovanni, 163</p> <p>Comba di Susa, 119</p> <p>Comfort as a moral influence, 185</p> <p>Comic song, the landlord's, 128</p> <p>Common sense, the safest guide, 108</p> <p>Consistent, who ever is? 153</p> <p>Contradictory principles, there must be a harmonious fusing of, 152</p> <p>Converting things by eating them, 153</p> <p>Corpses, desiccated, at S. Michele, 97</p> <p>Cousins, my, the lower animals, 69</p> <p>Cows fighting in farmyard, 120</p> <p>Crioco, 123</p> <p>Cristoforo, S., church of, at Mesocco, 208</p> <p>— — — at Castello, 234</p> <p>Crossing, efficacy of, 152</p> <p>— unexpected results of, 55</p> <p>— useless if too wide, 157</p> <p>Crucifixion, fresco at Fusio, 140</p> <p>Culture and priggishness, 141</p> <p>— a mode of concealing weakness, 192</p> <p>Current feeling, the safest guide, 108</p> <p>Cutlets, burnt, and the waiter, 141</p> <p>Dalpe, 38</p> <p>Dante a humbug, 156</p> <p>Darwin, Charles, no place for meeting, 69</p> <p>Darwin, Erasmus, 23, 131</p> <p>Dazio, Signor Pietro, of Fusio, 177</p> <p>Death, no man can die to himself, 277</p> <p>Deceit a necessary alloy of truth, 180</p> <p>Deportment, good technique resembles, 156</p> <p>Desire and power, 108</p> <p>Development of power to know our own likes and dislikes, 22</p> <p>Devil's Bridge, 23</p> <p>Diatonic scale, and song of birds in New Zealand, 232</p> <p>Dirt, eating a peck of moral, 71</p> <p><i>Disgrazia</i> and misfortune, 58</p> <p>D'Israeli, Isaac, quotations from, 67</p> <p>Dissenters all narrow-minded, 153</p> <p>Distribution of plants and animals often inexplicable, 135</p> <p>Diversion of mental images, 54</p> <p>Doera, fresco at, 145, 221</p> <p>Dogs, 156, 202, 260, 313</p> <p>Doing, the only mode of learning, 151</p> <p>Doors, how they open in time, 151</p> <p>Doubt, "There lives more doubt in honest faith," 67</p> <p>Downs, the South, like Monte Generoso, 230</p> <p>Draughtsman, first business of a, 148</p> <p>Drawing, the old manner of teaching, 150</p> <p>Dream, my, at Lago di Cadagno, 82</p> <p>Drunkenness and imagination, 46</p> <p><i>Duagne</i>, 133</p> <p>Duso, Agostino, his fresco at Sta. Maria in Calanca, 225</p> <p>Earnestness, 142, 192</p> <p>Eating, a mode of bigotry, 153</p> <p>Echo at Graglia, 192</p> <p>Edelweiss, 291</p> <p>Electricity and Alpine roads, 60</p> <p>Elephant brays a third, 233</p> <p>"Elongated" honey, 293</p> <p>Embryonic stages, the artist</p> <p>English priests and Italian, 106</p> <p>— why introspective, 18</p> <p>Equilibrium only attainable at the cost of progress, 195</p> <p><i>Eritis</i>, a panic concerning, 204</p> <p>Eternal punishment, 111, 196</p> <p>Eusebius, St., 178</p> <p>Evolution and illusion, 43</p> <p>— essence of, consists in not shocking too much, 110</p> <p>Extreme, every, an absurdity, 153</p> <p>Faldo, 22</p> <p>Faith, doubt lives in honest, 67</p> <p>— more assured in the days of spiritual Saturnalia, 68</p> <p>— foundations of our system based on, 107, 277</p> <p>— and reason, 108</p> <p>— catholic, of protoplasm, 152</p> <p>— a mode of impudence, 283</p> <p>Falsehood turning to truth, 71</p> <p>Famine prices at Locarno, 276</p> <p>Feeling, current, the safest guide, 108</p> <p>Fertile, rich and poor rarely fertile <i>inter se</i>, 195</p> <p>Fires, how Italians manage their, 117</p> <p>Fishmonger choosing a Monte, 23</p> <p>Flats and sharps, a maze of metaphysical, 23</p> <p>Fleet Street, beauties of, 19</p> <p>Flowers, names of, 291</p> <p>Fossil-soul, 234</p> <p>Foundations of action lie deeper than reason, 107</p> <p>— of a durable system laid on faith, 277</p> <p>Francis, St., and Insurance Co.'s plate, 191</p> <p>Friction, which prevents the unduly rapid growth of inventions, 60</p>
---	---

- Die CortexDB ist eine multi-model Datenbank, die für die Ablage von Datensätzen das Konzept eines document store nutzt, alle Felder und Feldinhalte redundanzfrei in einen mehrdimensionalen key/value-store überträgt (6. Normalform / 6NF), sowie Transaktionszeitpunkt und Gültigkeitszeitpunkt von Informationen bitemporal speichern kann. Darüber hinaus ist die ein-eindeutige Verbindung zwischen Datensätzen als Verweisstruktur möglich ("GraphDB").

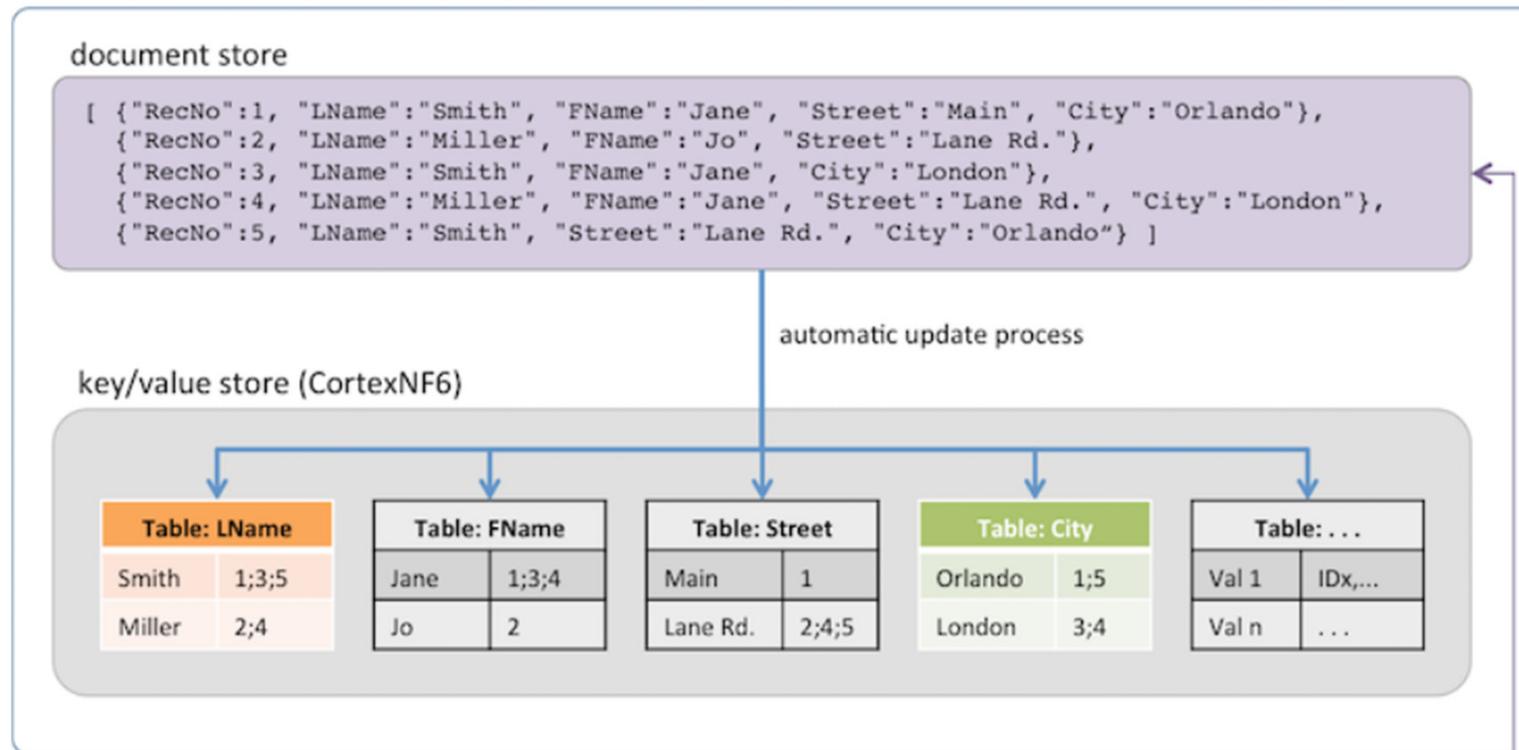
?

# CortexDB

cortex



CortexDB



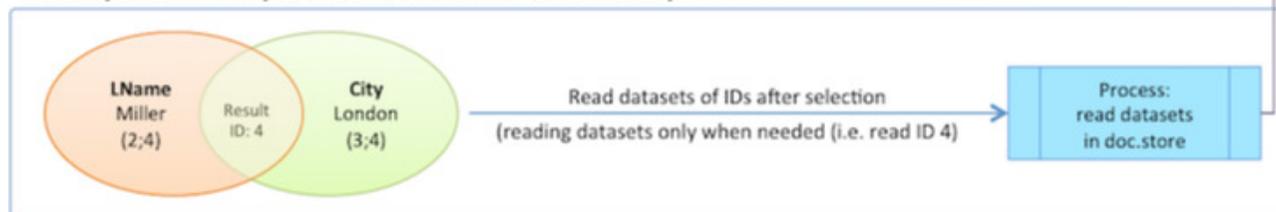
Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# CortexDB und Abfrage

cortex



developer selection process with methods of set theory



Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# CortexDB im Speziellen

cortex



- Bei der CortexDB handelt es sich um eine schemalose Datenbank (vergleichbar mit *Document Stores*), die mit einer inhaltsbasierenden mehrstufigen Indexstruktur arbeitet. Das bedeutet, dass jeder Inhalt *weiß*, in welchen Datensätzen und Feldern er vorkommt und jedes Feld *weiß*, welche verschiedenen Inhalte existieren.

# CortexDB im Speziellen

cortex



- Diese Indexstruktur bildet daher über alle Felder und deren Inhalte die gesamte Datenbank ab, so dass alle Abfragen in der Indexstruktur behandelt und die schemalos gespeicherten Daten nur zur Ausgabe verwendet werden. Die geringe Größe ermöglicht sehr schnelle Abfragen in beliebiger Abfragekombinationen. Administratoren und Entwickler können daher sofort mit diesem Index arbeiten und brauchen daher (auch bei Änderungen am Datenmodell) keinerlei Anpassungen oder Optimierungen vornehmen.

Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# CortexDB und Multi-Model



- Die CortexDB bietet die Möglichkeit, unterschiedliche Datenbank-Funktionen im Rahmen der schemalosen Struktur zu verwenden.
- So ist es beispielsweise möglich, dass über Verweisfelder Strukturen aufgebaut werden, wie sie bei Graph-Datenbanken zu finden sind. Weiterhin können einzelne Felder (optional) beliebig häufig in einem Datensatz verwendet werden (z.B. für Bankverbindungen oder eMail-Adressen); genauso können vergangene, aktuelle und künftige Werte je Feld eines Datensatzes gespeichert werden (bitemporale Datenbank). Zudem wird intern die 6. Normalform ähnlich in Anlehnung an Key/Value-Stores gespeichert (mehrdimensionaler Key/Value-Store).

Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# CortexDB und Multi-Model cortex



person

first name → John

last name → Doe

birth date → 27.07.1965

gender → male

links

1. → Jane, Doe ↔ wife

valid from	content		
27.07.1993	mutual friend		
03.09.1993	friend		
31.12.1993	girlfriend		
01.07.1994	fiancée		
03.09.1994	wife		

Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# CortexDB und Bitemporal

cortex



- Zu einem Datensatz speichert die CortexDB die Transaktionszeit von Änderungen. Ergänzend dazu kann jeder Inhalt aller Felder einzeln innerhalb eines Datensatzes mit einem Gültigkeitszeitpunkt versehen werden (Attribut-Zeitstempelung). Daher ist zu jedem Datensatz nicht nur der Zustand eines Datensatzes bei der letzten Transaktion ersichtlich, sondern auch der komplette "Lebenszyklus" (Historisierung). In diesem Fall spricht man von bitemporaler Datenbank.
- Dabei gilt, dass die Transaktionszeit automatisch (implizit) vom Server gepflegt wird und die Gültigkeit eines Wertes vom Anwender (explizit) bzw. von Automatismen wie dem Datenimport oder Schnittstellen. Wurde keine Gültigkeit gepflegt, gilt "unbekannt" für diesen Wert so lange, bis dieser mit einem neuen Wert überschrieben (eine Transaktion überschreibt diesen Wert) oder ein weiterer Wert mit einer Gültigkeit hinzugefügt wird.

Cortex AG <https://docs.cortex-ag.com/de/> heruntergeladen am 24.10.2018

# Identity Mapping / Record Linkage / Dynamic Links in Data Vault



- Mögliche dynamische Verbindungen werden von Algorithmen bewertet und von Maschine Learning Algorithmen bestimmt.

## Das Problem

FIRMA DATENQUELLE 1	FIRMA DATENQUELLE 2
Ähnlichkeit Algorithmus GmbH	Ähnlichkeit Algorythmus GmbH
Completely Different Limited	Völlig Unterschiedlich Limited
Der ganz andere Verein eV	Der völlig andere Verein e.V.

# Kein gemeinsamer Identifikator



- Es gibt mehrere aktualisierbare Stammdatenquellen (z.B. Firmendaten).
- Es gibt keinen gemeinsamen Schlüssel für einfache Joins.
- Die Datenzusammenführung muss auf der Grundlage der Zeichenkettenähnlichkeit erfolgen.
- - Name, Adresse, E-Mail, Ort, etc.
- Welcher Ähnlichkeitsalgorithmus ist zu verwenden?
- Identitätsabbildung und Deduplizierung stellen tatsächlich das gleiche Problem dar.

# Trade-Offs



- Wie viele Daten möchten Sie automatisch abgleichen?
- Je mehr Daten Sie automatisch abgleichen, desto mehr Fehler können Sie erwarten.
- Wie viele Fehler können Sie sich mit dem automatischen Algorithmus leisten?
- Unmöglich, automatisch 100% korrekte Übereinstimmungen zu erhalten.
- Für eine Genauigkeit von 100% müssen die Data Stewards den endgültigen Abgleich manuell durchführen.

# Performance Probleme



- Für ein ungefähre Zusammenführen kann jede Zeile aus einer Tabelle mit jeder Zeile aus einer anderen Tabelle verbunden werden.
- Cross Join?
- Selbst bei kleinen Datensätzen können Cross Joins zu Performance-Problemen führen.
- Daten müssen in kleineren Teilen abgebildet werden.
- Partitionierung
- Sortieren der Nachbarschaft
- Datenmenge verringern

# Iterativer Ansatz



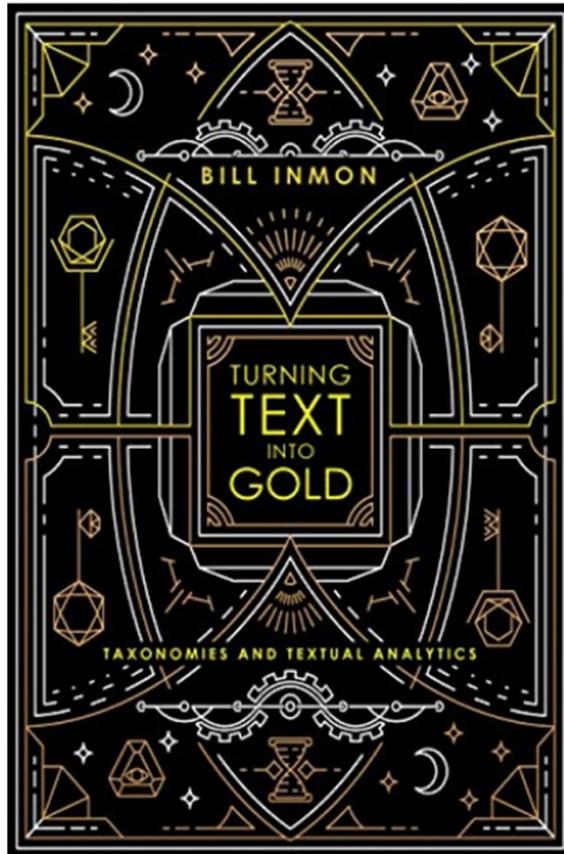
- Matching in Iterationen durchführen
- Beginnen Sie mit exakten Übereinstimmungen (Inner Join)
- Dann fügen Sie ähnliche Zeichenketten zusammen.
- Dann fügen Sie weniger ähnliche Zeichenketten zusammen.
- Dann manuell zusammenführen
- Die Lösung hängt vom Werkzeug ab.
- Bereinigen Sie die Daten so weit wie möglich, bevor Sie sie abgleichen. ???

# MDS Ähnlichkeitsalgorithmen



- Master Data Services implementiert einige bekannte Ähnlichkeitsalgorithmen.
- Levenshtein-Distanz (auch bekannt als Edit-Distanz)
- Jaccard-Index
- Jaro-Winkler Abstand
- Simil-Algorithmus (auch bekannt als Ratcliff/Obershelp)
- NGrams
- Implementiert in `mdq.NGrams` und `mdq.Similarity CLR` Funktionen

# Empfehlung



- Es kann ein komplexer Prozess sein, um Text in ein strukturiertes Format zu konvertieren. Die Hauptschwierigkeit besteht darin, mit Unklarheiten im Zusammenhang mit dem Text umzugehen, was erfordert, dass wir verschiedene Techniken wie Homographieauflösung, Näherungsanalyse, Negativitätsableitung usw. anwenden.

Uli Bethke <https://sonra.io/2018/10/11/the-future-of-etl-and-the-limitations-of-data-virtualisation-and-noetl/>



**Vielen Dank für ihre  
Aufmerksamkeit**

