

# Data Engineer

**Data-Informationen:**

Data Governance, data Vault, ..

## Definition von Data Governance

Data Governance umfasst alle Verfahren, mit denen für die Sicherheit, Korrektheit, Verfügbarkeit und Nutzbarkeit von Daten sowie den Datenschutz gesorgt wird. Dazu gehören die Maßnahmen, die Mitarbeiter ergreifen müssen, die Prozesse, die sie befolgen müssen, und die Technologien, die sie während des Datenlebenszyklus unterstützen.

<https://www.kobold.ai/data-governance/>

<https://barc.com/de/wege-zur-data-governance-teil-1/>

<https://barc.com/de/wege-zur-data-governance-teil-2/>

<https://barc.com/de/wege-zur-data-governance-teil-3/>

## Welche Vorteile bietet Data Governance?

1. Bessere und schnellere Entscheidungen
2. Bessere Kostenkontrolle
3. Bessere Einhaltung gesetzlicher Vorschriften
4. Gestärktes Vertrauen bei Kunden und Lieferanten
5. Optimiertes Risikomanagement
6. Mehr Mitarbeiter können auf mehr Daten zugreifen.

<https://www.ibsolution.com/academy/blog/master-data-management/praxisbeispiele-fuer-den-nutzen-von-data-governance>

# Data Governance 3



## Fragen werden zu Lösungen durch Data Governance

Bei der Datenstrategie-Entwicklung suchen Unternehmen häufig Antworten auf beispielsweise nachfolgende Fragen. Diese und weitere Punkte regelt Data Governance:

Auf welcher Basis definiere ich meine (datenbasierende) Anforderung? Was gibt es schon? Wer ist autorisiert das zu verändern?	Was bedeutet diese Kennzahl? Kann ich diese Daten für meine Analyse/Entscheidung verwenden? Wer weiß, was das bedeutet?
Wir wissen nicht genau wo bestimmte Daten herkommen, wie sie kombiniert werden und können dadurch Datenbestände nicht nutzen.	Wir haben ein neues Data Warehouse und können es nicht sinnvoll weiterentwickeln, da wir nicht wissen worauf zu achten ist.
Wir haben das Vertrauen in die Daten des Data Warehouse verloren, da einige Zubauten von unterschiedlichen Personen/Firmen durchgeführt wurden.	Die Datenqualität in unserem Data Warehouse ist nicht transparent, wir kennen Sie nicht.
Wer kennt sich fachlich mit den Daten aus? Wer legt fest was gute Daten sind und wer ist verantwortlich für die Qualität der Daten?	Was ist rechtlich alles zu beachten? Und dürfen die Daten für Analysen verwendet werden?
Woher kommen Daten für Kennzahlen und wie wird sie errechnet? (Lineage)	Woher bekomme ich die Daten? Und wer genehmigt den Zugriff darauf?

## **Was Data Governance NICHT ist!**

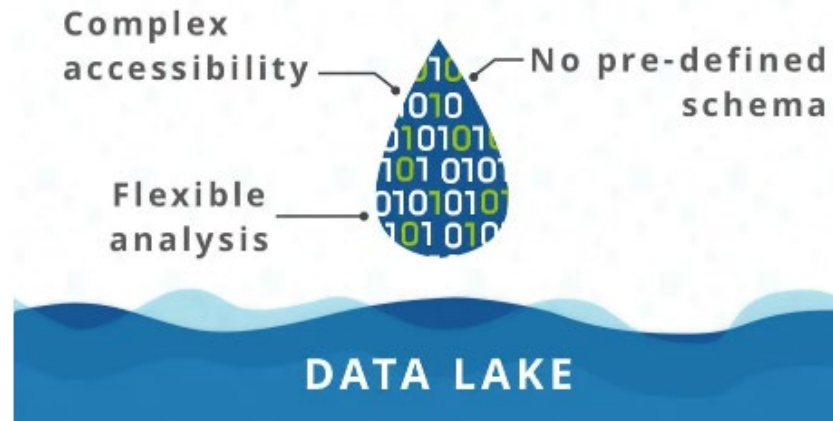
1. Data Governance ist nicht Datenmanagement
2. Data Governance ist nicht Stammdatenmanagement
3. Data Governance ist nicht Data Stewardship

## Was Data Warehouses nicht sind.

1. Data Warehouses sind keine Datenbanken
2. Data Warehouses sind keine Data Lakes

	Datenbank	Data Warehouse
Definition	Enthält gesammelte Daten für unterschiedliche Transaktionszwecke. Ist auf den Lese- und Schreibzugriff optimiert.	Enthält aggregierte Daten, die für Analysezwecke transformiert werden. Ist auf die Aggregation (Verdichtung) und den Zugriff auf große Datensätze optimiert.
Verwendung	Datenbanken sind darauf ausgelegt, Informationen schnell zu erfassen und abzurufen.	Data Warehouses speichern Daten aus mehreren Datenbanken, um die Analyse zu vereinfachen.
Arten	Es gibt verschiedene Arten von Datenbanken, z. B. csv-, html- und Excel-Tabellen. Transaktionale Online-Verarbeitungsdatenbanken kommen auch beim Data Warehousing zum Einsatz.	Data Warehouses bestehen aus analytischen Datenbanken. Diese bauen wiederum Transaktionsdatenbanken auf, um Analysen zu ermöglichen.

# Data Lake 1



## Data Lake Definition:

Ein Data Lake ist ein zentrales Repository, das Big Data aus unterschiedlichen Quellen in einem rohen, granularen Format speichert. Es kann strukturierte, semistrukturierte oder unstrukturierte Daten aufnehmen. D. h. die Daten können in einem flexibleren Format zur späteren Nutzung aufbewahrt werden. Ein Data Lake verbindet beim Speichervorgang Daten mit Identifiern und Metadaten-Tags, um einen schnelleren Zugriff zu gewährleisten.

[Video zu Data Lake](#)

<https://www.bigdata-insider.de/was-ist-ein-data-lake-a-686778/>

# Data Warehouse vs. Data Lake

	Data Lake	Data Warehouse
Datenstruktur	Roh	Verarbeitet
Zweck der Daten	Noch nicht festgelegt	Aktuell in Gebrauch
Benutzer	Data Scientists	Business-Anwender
Zugänglichkeit	Gut zugänglich und schnell zu aktualisieren	Komplizierter und teuer, Änderungen vorzunehmen

<https://www.kobold.ai/data-warehouse-vs-lake/>

<https://www.weclapp.com/de/lexikon/data-warehouse/>

<https://www.bigdata-insider.de/vom-data-warehouse-und-data-lake-zum-lakehouse-a-834c419b7019aeb513bfad691fb37202/>



## Was ist Business Intelligence?

Der Begriff „Business Intelligence“ beschreibt letztendlich einen Prozess, nämlich die Erfassung, Speicherung und Analyse von Daten aus geschäftlichen Vorgängen.

BI bietet dabei umfassende geschäftliche Kennzahlen in beinahe Echtzeit zur Unterstützung besserer Entscheidungen.

Sie können mit besserer Business Intelligence Leistungs-Benchmarks erstellen, Markttrends ermitteln, die Compliance erhöhen und beinahe jeden Aspekt Ihres Unternehmens optimieren.

## Was ist Data Analytics?

Data Analytics bezeichnet den technischen Vorgang zur Ermittlung von Daten, zur Aufbereitung von Daten, zur Transformation von Daten und zum Aufbau eines Systems zur Verwaltung von Daten. Mit Data Analytics wird versucht, aus großen Datenmengen Trends zu ermitteln und Probleme zu lösen. Data Analytics wird in vielen Bereichen verwendet – in Behörden ebenso wie in der Wissenschaft. Diese Technik ist also nicht auf geschäftliche Anwendungen beschränkt.

## Was ist Data Science?

Data Science ist die Fähigkeit einer Organisation, Datenmengen zu analysieren, zu extrahieren und zu formatieren, um sie auf visuelle und aussagekräftige Weise zu präsentieren.

Es geht darum, zukunftsorientierte Trends zu erheben und darzustellen.

Ein Data Scientist übernimmt diese Aufgabe, anhand von aussagekräftigen Daten Wege zur Beantwortung von Fragen auszuarbeiten und sich Hypothesen für die Zukunft auszudenken.

Die Data Science gilt als reaktiv und ist in der Antizipation angesiedelt.

# Descriptive, Diagnostic, Predictive und Prescriptive Analytics



Die **Descriptive Analytics** versucht zu „beschreiben“, wie sich der Umsatz in der Vergangenheit entwickelt hat.

Die **Diagnostic Analytics** versucht zu „diagnostizieren“, warum sich der Umsatz positiv oder negativ entwickelt hat.

Die **Predictive Analytics** versucht „vorherzusagen“, wie sich der Umsatz in den nächsten Monaten oder Jahren entwickeln wird.

Die **Prescriptive Analytics** versucht „vorschreiben“, wie sichergestellt werden kann, dass sich der Umsatz weiter positiv und nicht negativ entwickelt.

# Abgrenzung der Begriffe zueinander



[Business Intelligence vs. Data Science](#)

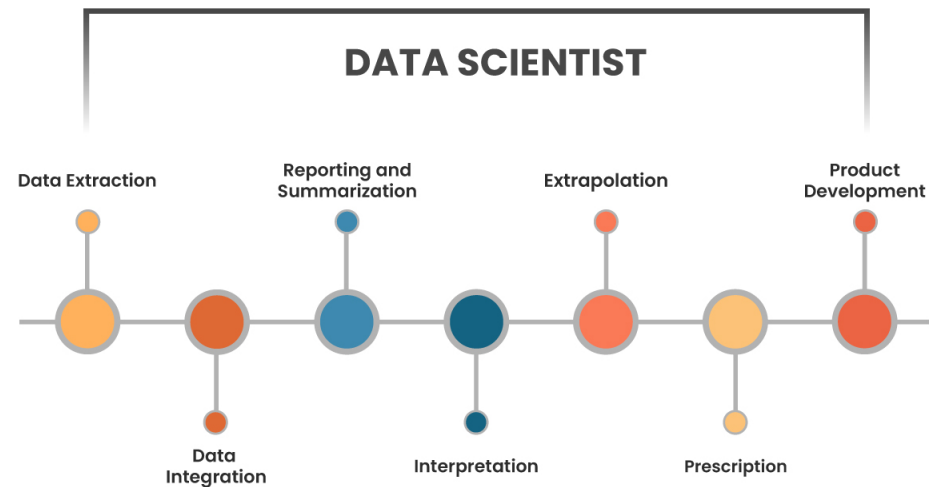
[Data Science vs. Business Intelligence](#)

[Business Analytics vs. Data Science](#)

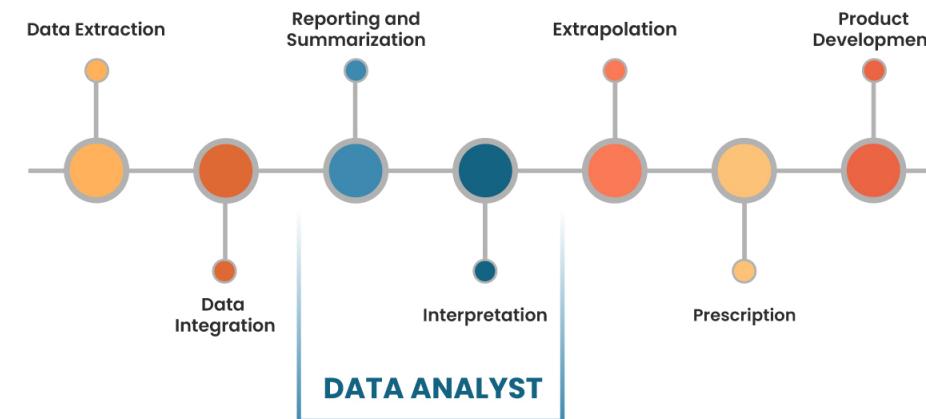
[Business Intelligence vs. Business Analytics vs. Big Data Analytics](#)

[Data Analytics vs. Data Science](#)

# Data Analytics vs. Data Science 1



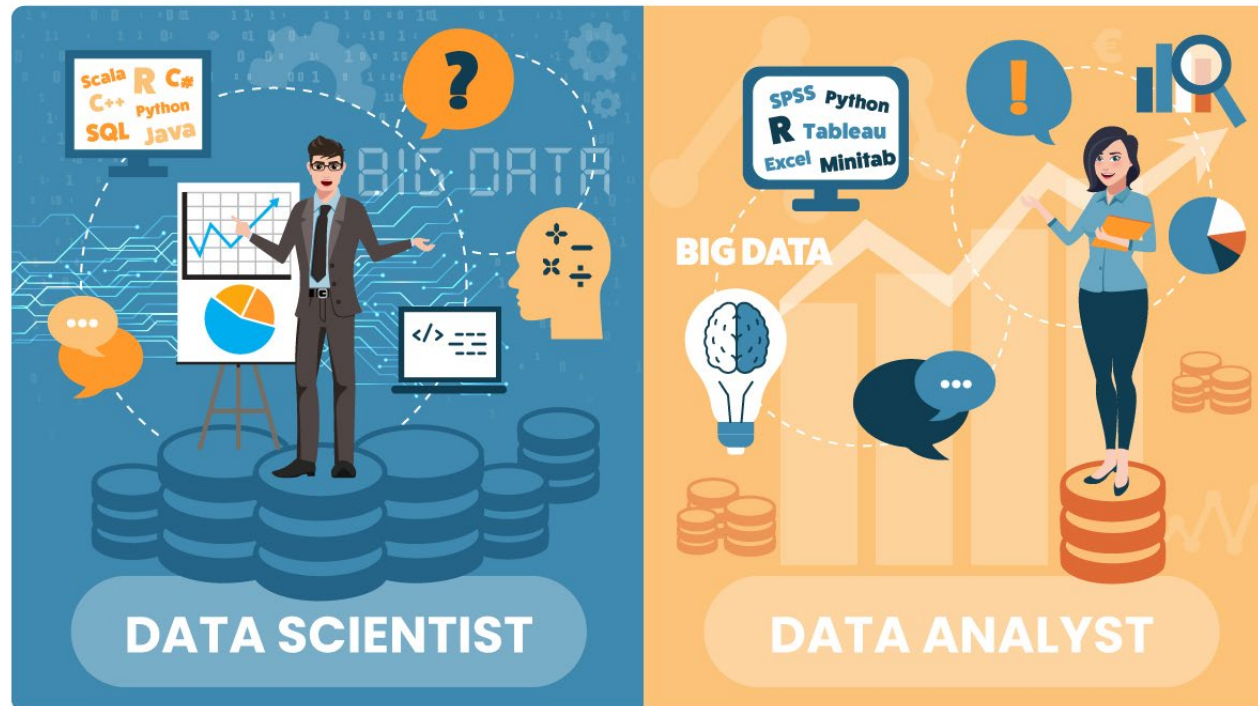
Wirtschaftsforum in Kooperation mit Westphalia DataLab



Wirtschaftsforum in Kooperation mit Westphalia DataLab

## Data Analytics vs. Data Science

# Data Analytics vs. Data Science 2



# Business Intelligence und Data Science – Hand in Hand



Auch wenn die Data Science dank ihrer Fähigkeit zur Zukunftsforschung derzeit im Aufwind ist, verliert sie deutlich an Relevanz, wenn sie sich nicht auf die von der BI gelieferten Analysen stützt.

Wie Victor Hugo schon sagte: „Die Zukunft ist eine Tür, die Vergangenheit ist der Schlüssel dazu“. Mit anderen Worten: Die BI muss mehr denn je die Grundlage für die Datenwissenschaft sein. Die Data Science kann sich dann auf das Bestehende stützen, um ihre Hypothesen zu untermauern.

Damit sich die Mitarbeiter einer Organisation selbstständig mit BI beschäftigen können, gibt es Lösungen wie Power BI von Microsoft. Die Nutzung kann jedes Unternehmen selbst durchführen, da man kein IT-Experte sein muss, um mit dieser Lösung Big-Data-Daten zu verarbeiten.

Die von der Datenwissenschaft angebotenen Tools werden den verschiedenen hierarchischen Ebenen eines Unternehmens Ratschläge geben, wie sie dieses Wissen am besten nutzen können.



# SQL vs. NoSQL

	SQL	NoSQL
Beschreibung	relational	nicht-relational
Anwendung	Abfrage zum Analysieren und Abrufen von Daten	für eine Vielzahl moderner Anwendungen wie WebApps geeignet
Abfragesprache	SQL	mehrere Sprachen je nach Anwendung
Typ	Tabelle	Dokument / Graph / Key-Value
Schema	festgelegt und vordefiniert	dynamisch
DMS(Beispiele)	Oracle, PostGres, MySQL	MongoDB, Neo4J
Eignet sich für	komplexe und intensive Abfragen	Große Datenbanken, Big Data
Entwicklungsjahre	70er Jahre	2000er
Open Source	Open Source (PostGres, MySQL) und proprietäre Systeme (Oracle)	Open Source
Vorteile	optimierte Datenspeicherung und Stabilität	einfache und flexible Speicherung
Nachteile	keine Flexibilität, erforderliche Expertise	manchmal zu flexibel

# SQL oder NoSQL – Was sollte ich wählen?



## **SQL:**

**Du willst eine strukturierte und segmentierte Datenbank** (Grundlage der relationalen Datenbanken).

**Typ und Gültigkeit der Daten sind sehr wichtig**

**Du willst bestimmte Datenelemente schreiben und ändern**

(mit SQL kannst Du bestimmte Zeilen einfach ändern).

**Du brauchst komplexe Abfragen**

## **NoSql:**

**Du willst eine Datenbank ohne spezifisches Schema**

(Eine nicht fixierte Struktur zum Beispiel).

**Du brauchst viele Leseabfragen.** Mit NoSQL können tatsächlich alle benötigten Daten auf einmal abgerufen werden, ohne besondere Joins.

**Große Datensätze** (Big Data)

**Verteilte Daten** (mehrere Quellen)

Cloud-Computing, häufig auch einfach als „die Cloud“ bezeichnet, ist die Bereitstellung von IT-Ressourcen – von Anwendungen bis zu Rechenzentren – bei Bedarf über das Internet auf der Basis nutzungsabhängiger Gebühren.

- Flexible, skalierbare Ressourcen, die schnell und einfach an den Bedarf angepasst werden können
- Nutzungsabhängiger Service, bei dem Sie nur für die tatsächlich genutzten Ressourcen bezahlen
- Self-Service-Zugriff auf alle IT-Ressourcen, die Sie benötigen

<https://www.timetac.com/de/zeiterfassungslexikon/cloud-computing/>

## Cloud Services sind:

- Software as a Service (SaaS)
- Plattform as a Service (PaaS)
- Infrastructure as a Service (IaaS)
- Public Cloud
- Private Cloud
- Hybrid Cloud

<https://www.ibm.com/de-de/cloud/learn/cloud-computing-gbl>

## Cloud Computing — Vorteile

- Kein Vorhalten der IT-Infrastruktur (SW u. HW)
- Sehr gute Skalierungsmöglichkeiten
- Nicht mehr selbst aktualisieren
- Klares Lizenzmanagement

## Cloud Computing — Nachteile

- Defekter Internetzugang
- Hacker
- Verschlüsselung, wie sicher?

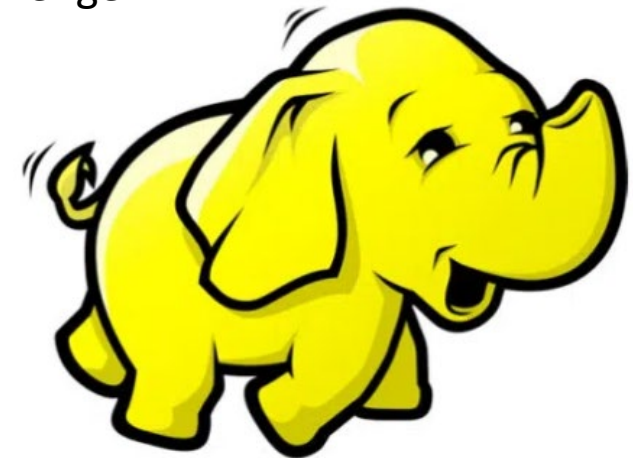
<https://www.geld-online-blog.de/software-tools/cloud-computing-begriffserklaerung-vor-und-nachteile/>

## Hadoop

Hadoop ist ein Java-basiertes Open-Source-Framework der Apache Software Foundation. Das Framework dient der Verarbeitung von großen Datenmengen, die auf verschiedene Systeme verteilt sind.

Zu den Stärken der Software gehört die hohe Geschwindigkeit der Big Data Datenverarbeitung. Aber auch aufgrund seiner intuitiven Oberfläche ist Hadoop eines der meistgenutzten Software Frameworks für die Verarbeitung großer Datenmengen.

<https://compamind.de/knowhow/hadoop/>



# Spark



Apache Spark ist eine einheitliche In-Memory Analytics Plattform für Big Data Verarbeitung, Data Streaming, SQL, Machine Learning und Graph Verarbeitung.

<https://datasolut.com/was-ist-spark/>



# Hadoop vs. Spark



<https://www.scnsoft.de/blog/spark-vs-hadoop-mapreduce>

<https://www.tecchannel.de/a/apache-spark-versus-hadoop,3197890>



# MapReduce 1

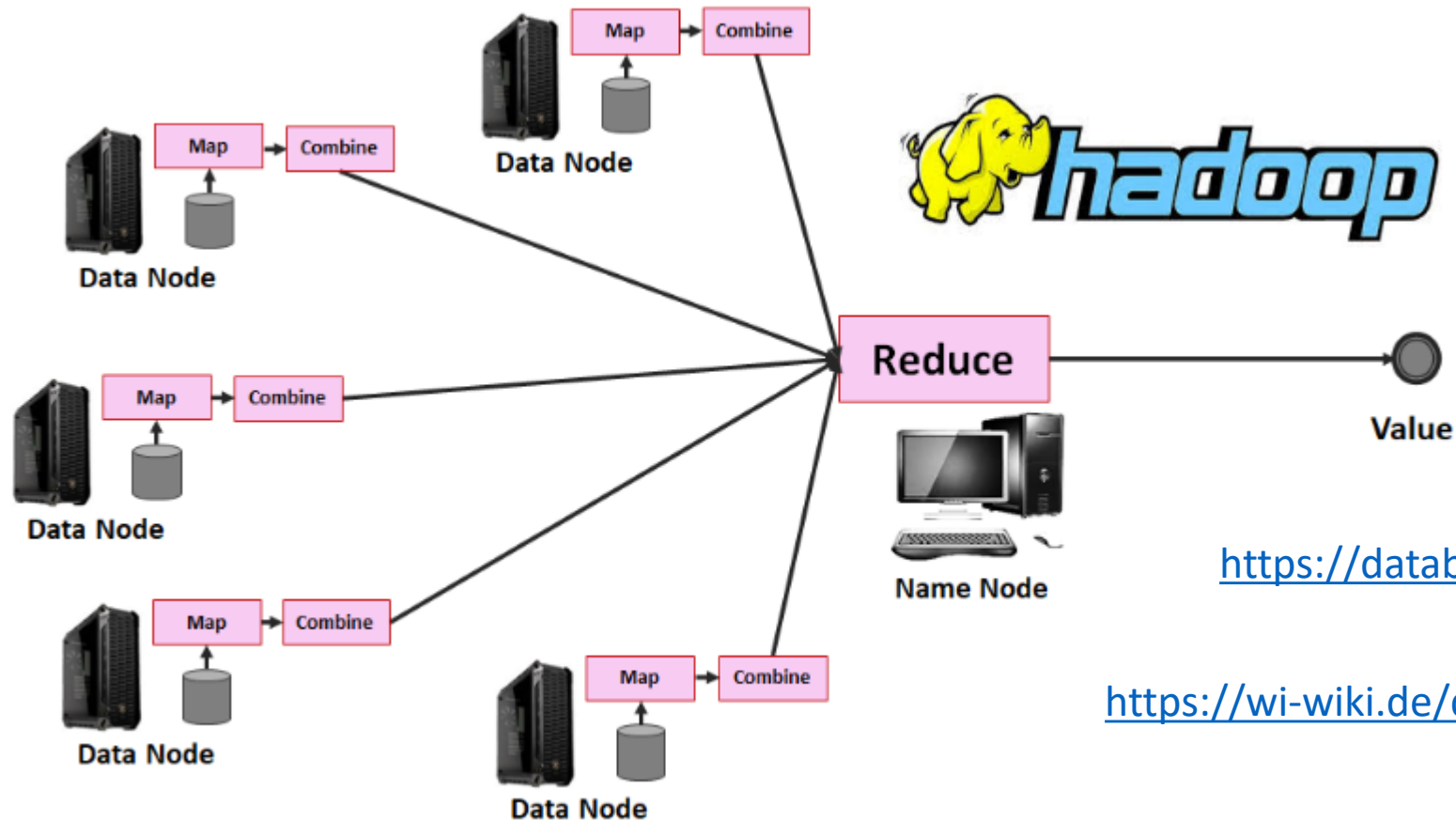


- Programmiermodell zur Verarbeitung von großen unstrukturierten oder semi-strukturierten Datensätzen
- MapReduce nutzt verteilte Speicherung der Daten in Blöcken (sorgt für Datenqualität bei fehlerhaften Schreib- oder Lesevorgängen)
- MapReduce-Framework sorgt für die Aufteilung der Berechnungen auf mehrere Recheneinheiten
- Dadurch parallele Ausführung auf mehreren Rechnern
- Nach Beendigung der Berechnungen aggregiert das Framework die Ergebnisse
- Entwickler von verteilten Anwendungen müssen nur das Framework benutzen, keine Codeänderungen bei der Änderung der Client-Anzahl nötig
- Verwendung von handelsüblichen Computern möglich
- Keine Notwendigkeit für spezielle High-End Server

[MapReduce.pdf](#)

[MapReduce Konzept.pdf](#)

# MapReduce 2

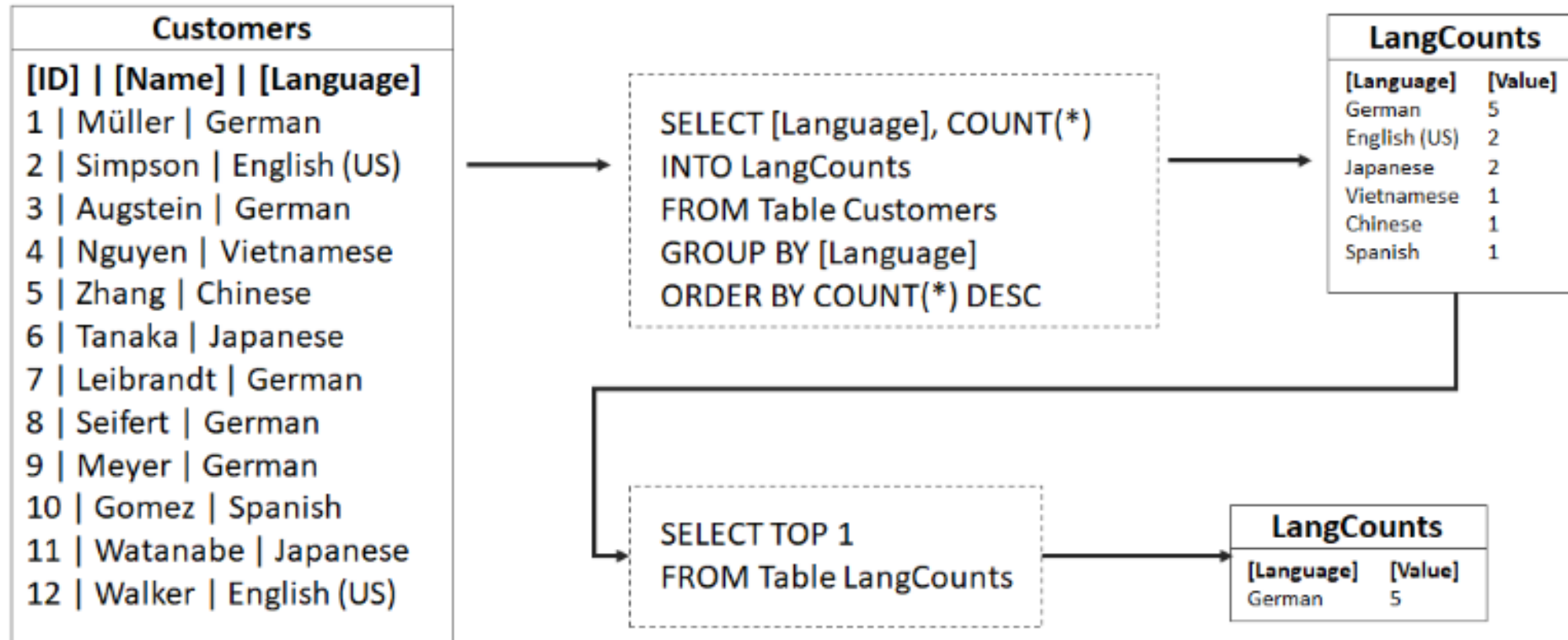


<https://databasecamp.de/daten/mapreduce>

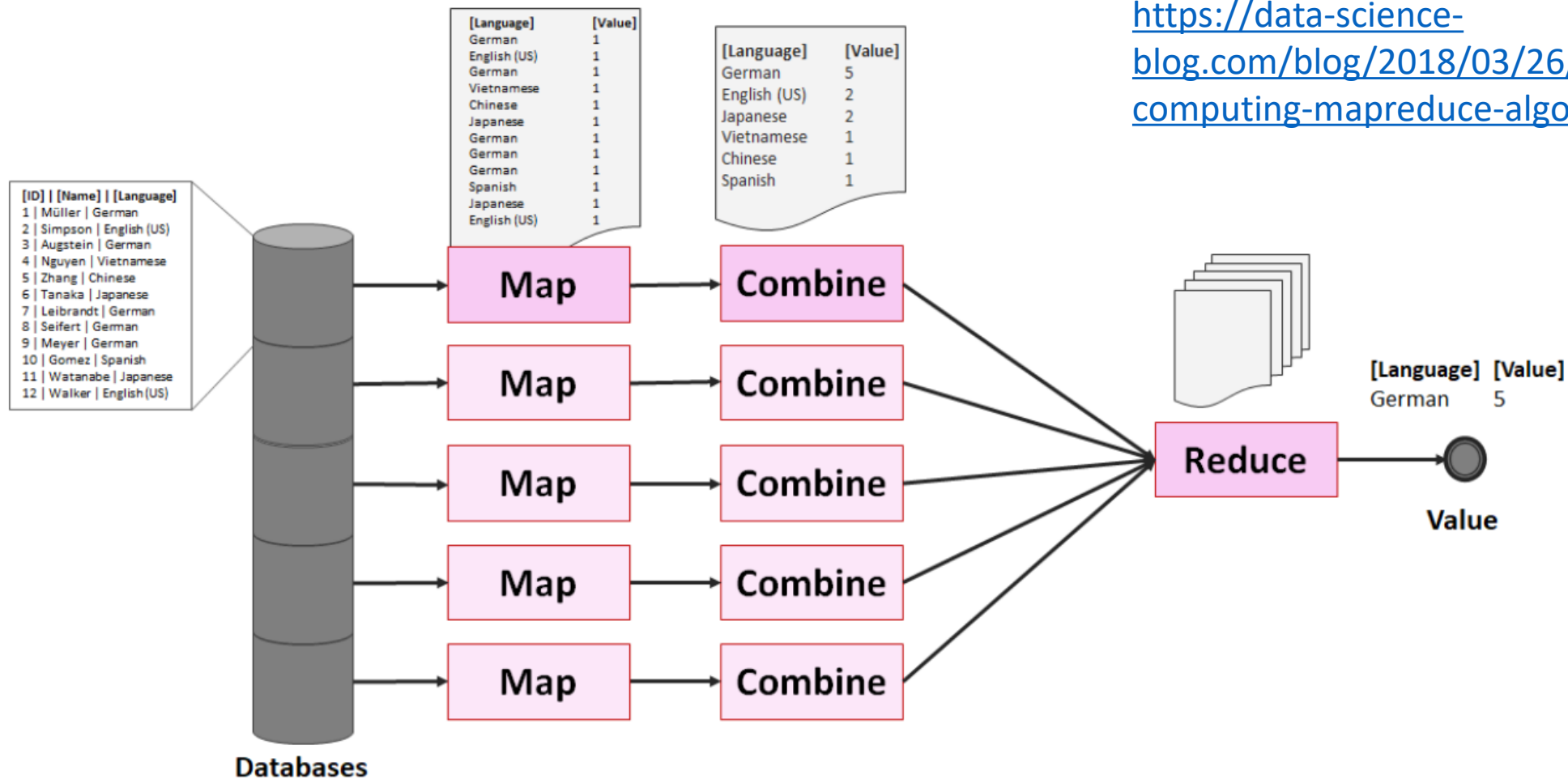
<https://wi-wiki.de/doku.php?id=bigdata:mapreduce>

<https://www.linux-magazin.de/ausgaben/2012/04/hadoop/>

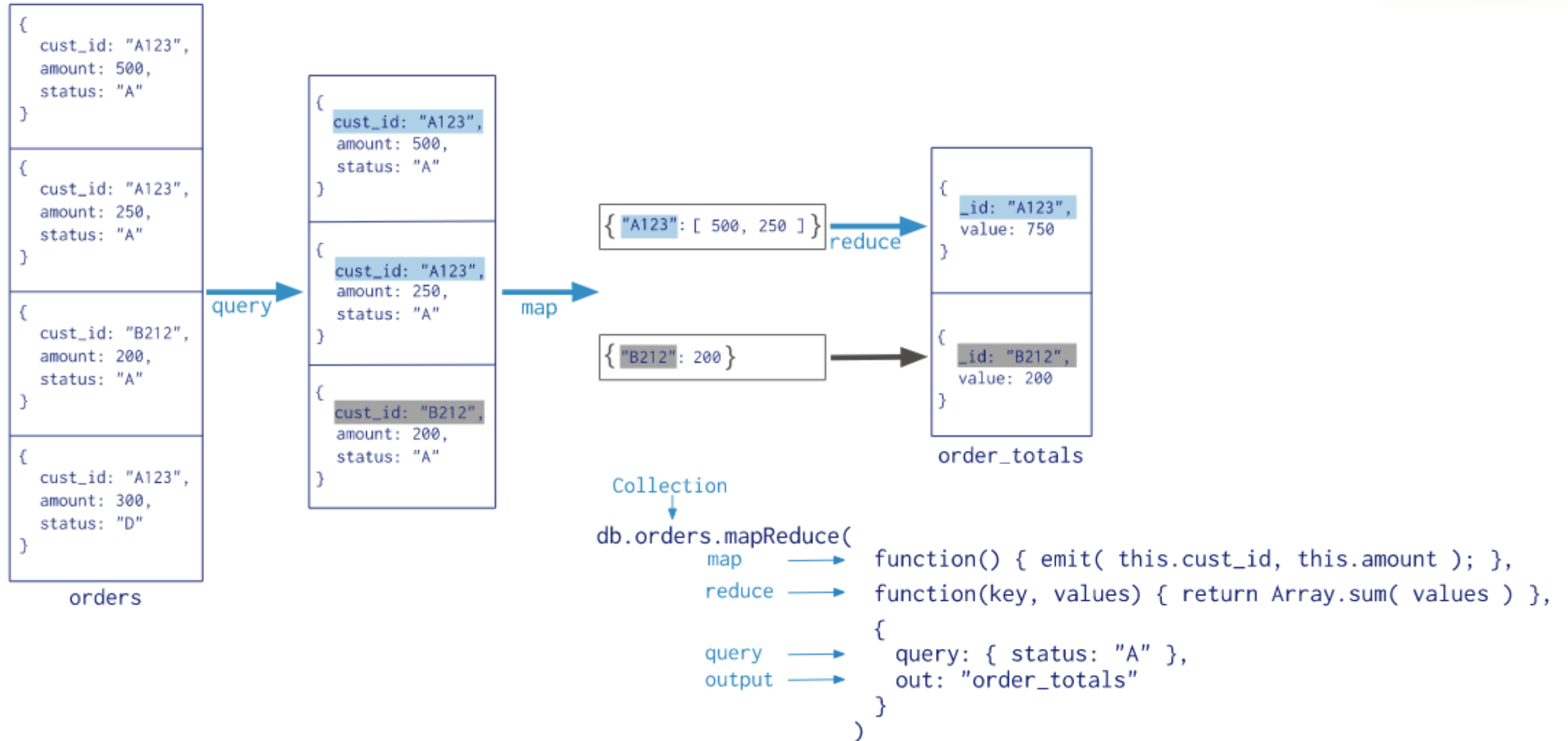
# MapReduce 3



# MapReduce 4



# MapReduce 5



Es gibt verschiedene Arten von Datenbanksystemen:

- **relationale Datenbanksysteme**
- **NoSQL-Datenbanksysteme**
- **grafbasierte Datenbanksysteme**
- **Objektdatenbanken**
- **In-Memory-Datenbanken**

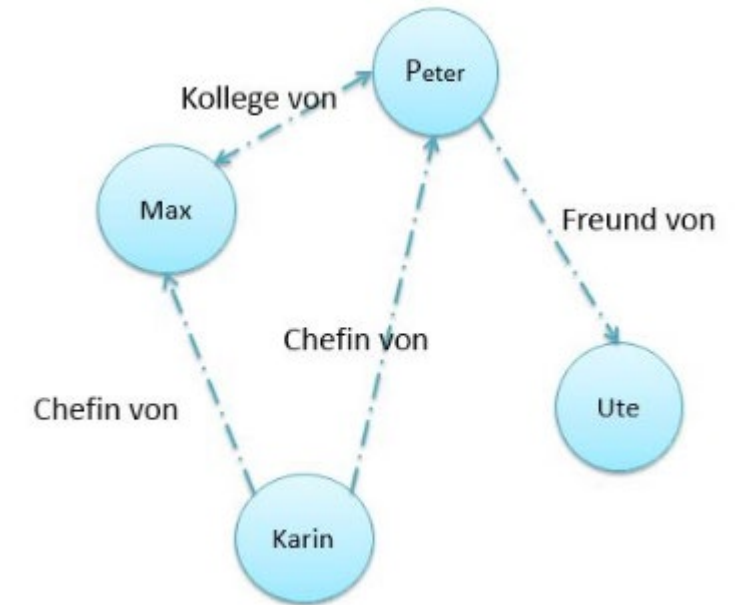
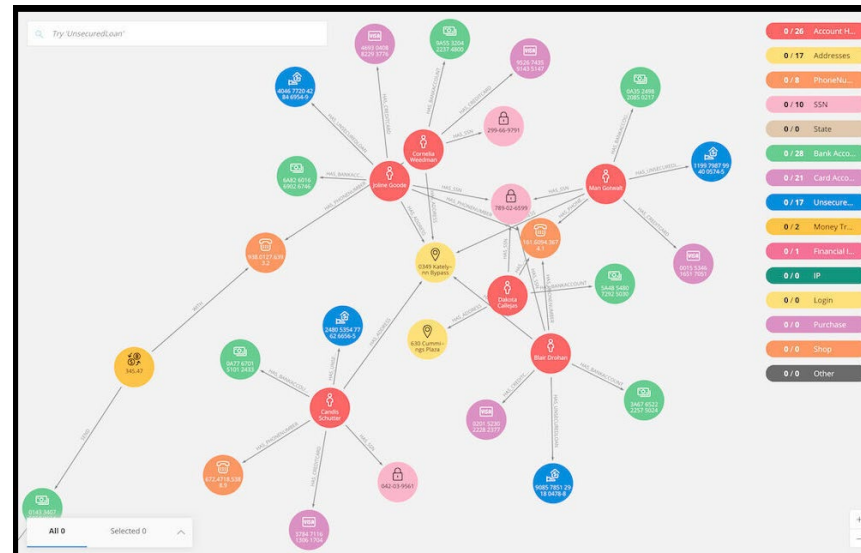
<https://www.ionos.de/digitalguide/hosting/hosting-technik/datenbanken/>

# Grafdatenbank

Eine Grafdatenbank ist eine Datenbank, die Graphen benutzt, um stark vernetzte Informationen darzustellen und abzuspeichern.

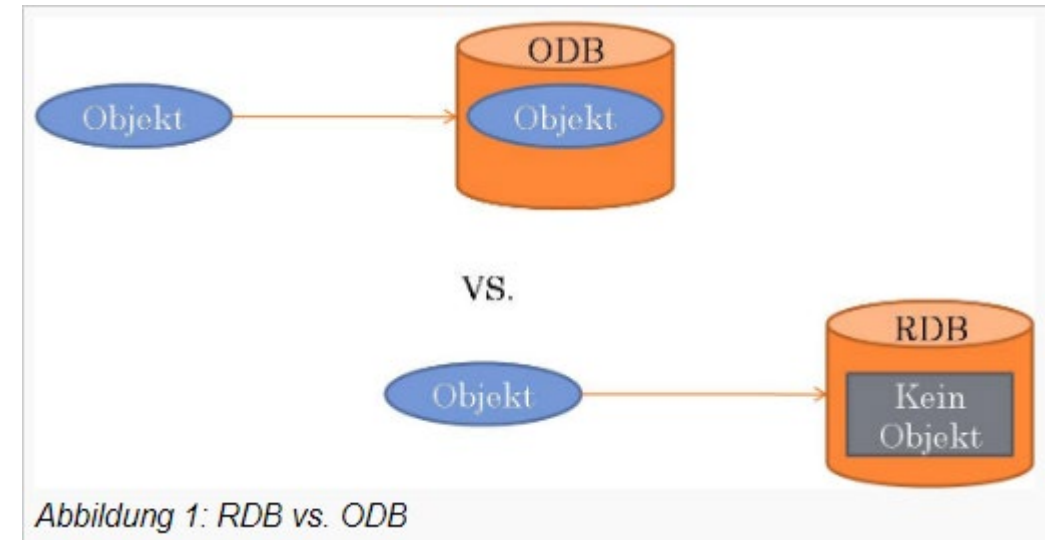
Ein solcher Graph besteht aus Knoten und Kanten, den Verbindungen zwischen den Knoten.

Die zwei bekanntesten Konzepte für Grafdatenbanken sind das Resource Description Framework (RDF) und Labeled-Property Graph (LPG).



Eine Objektdatenbank oder objektorientierte Datenbank ist eine Datenbank, die auf dem Objektdatenbankmodell basiert. Im Unterschied zur relationalen Datenbank werden Daten hier als Objekte im Sinne der Objektorientierung verwaltet. Das zugehörige Datenbankmanagementsystem wird als das objektorientierte Datenbankmanagementsystem bezeichnet. Objektdatenbank und Objektdatenbankmanagementsystem bilden gemeinsam das Objektdatenbanksystem.

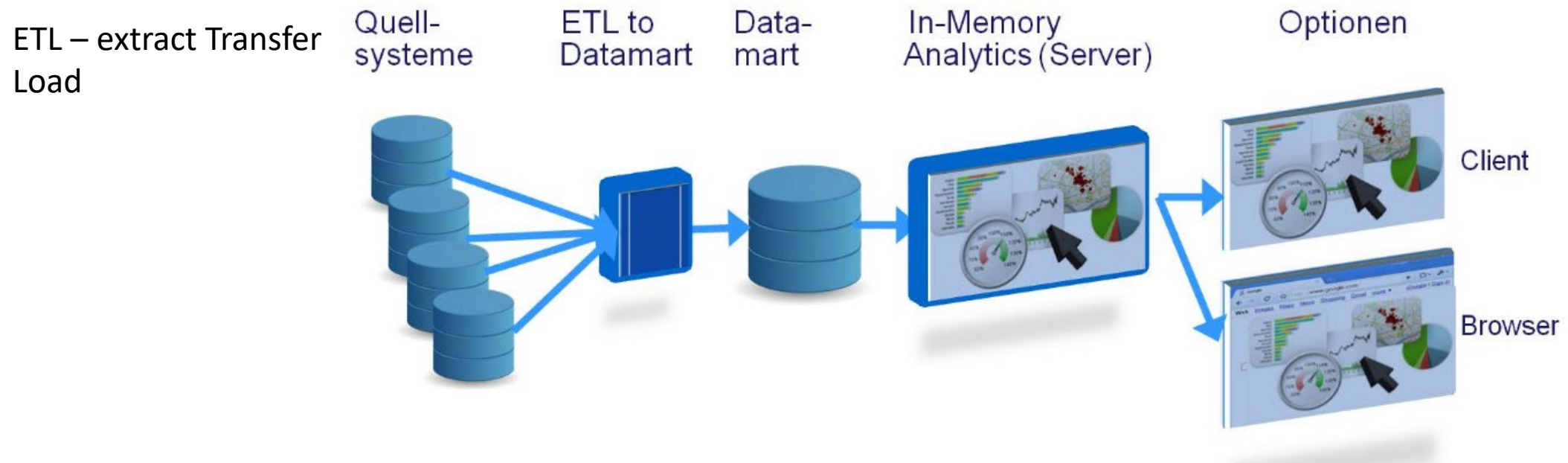
## Caché-Datenbanksystem





# In-Memory-Datenbank

Eine In-Memory-Datenbank (IMDB) ist ein Datenbankmanagementsystem, das den Arbeitsspeicher eines Computers als Datenspeicher nutzt. Damit unterscheidet es sich von herkömmlichen Datenbankmanagementsystemen, die dazu Festplattenlaufwerke verwenden.



Datentyp definiert Art – Zeichen, binär, numerisch usw ... Länge bzw. Größe Genauigkeit bei numerischen Daten Dezimalstellen bei numerischen Daten.

## **Kategorien von Datentypen:**

- Datentypen für Zeichenketten
- Numerische Datentypen
- Datentypen für Datums- und Uhrzeitwerte

# Basisdatentypen

Kategorie	Basisdatentyp	Beschreibung
Alphanumerisch	CHARACTER(n)	Jede Art von alphanumerischen Zeichen, die vom Zeichensatz bereitgestellt werden. Die Anzahl der Zeichen werden für die Zeichen in der Definition reserviert.
	CHARACTER VARYING(n)	Aufgebaut wie CHARACTER(n), reserviert der CHARACTER VARYING(n) aber nur die tatsächlich verwendeten Zeichen. Beispiel: Varchar(max) oder Varchar2
Numerisch	SMALLINT	Ganze Zahl (2 Byte)
	INTEGER	Ganze Zahl (4 Byte)
	BIGINT	Ganze Zahl (8 Byte)
	NUMERIC(m,n)	Dezimalzahl mit der festen Länge m und n Nachkommastellen
	FLOAT(n)	Gleitkommazahl mit gegebener Genauigkeit n
	REAL	Gleitkommazahl mit einfacher Genauigkeit
	DOUBLE	Gleitkommazahl mit doppelter Genauigkeit
Datum/Uhrzeit	DATE	Datumsangabe im Format „JJJJ-MM-TT“, abhängig von den Einstellungen
	TIME	Zeitangabe im Format „hh:mm:ss“, abhängig von den Einstellungen
	TIMESTAMP	Kombination aus DATE und TIME; wird oft für Zugriffsroutinen (Concurrency) benutzt
Intervall	INTERVAL HOUR	Intervall für Stunden, zum Beispiel alle 60 Minuten
	INTERVAL DAY	Intervall für Tage, zum Beispiel jeden 7 Tag
Sonstige	BLOB	Binary Large Object, dient zur Speicherung binärer Daten, wie Bilder, Tonaufnahmen oder ganzen Filmen
	BIT	Bitkette der Länge n
	BOOLEAN	Wahrheitswert für Wahr/Falsch

# Spezielle Datentypen



Spezielle Datentypen in Datenbanken sind Datentypen,  
die speziell an die Bedürfnisse einer Datenbank angepasst sind.

Dazu gehören beispielsweise Datentypen wie Character Large Object (CLOB), Timestamps [1],  
JSON Daten und Binärdaten.

CLOBs werden zur Speicherung von sehr langen Zeichenketten verwendet,

Timestamps für datenbankweit eindeutige Ausdrücke,

JSON Daten werden häufig in NoSQL oder anderen speziellen Datenbanken gespeichert und

Binärdaten werden zur Verschlüsselung von Daten in Datenbanken verwendet [2]

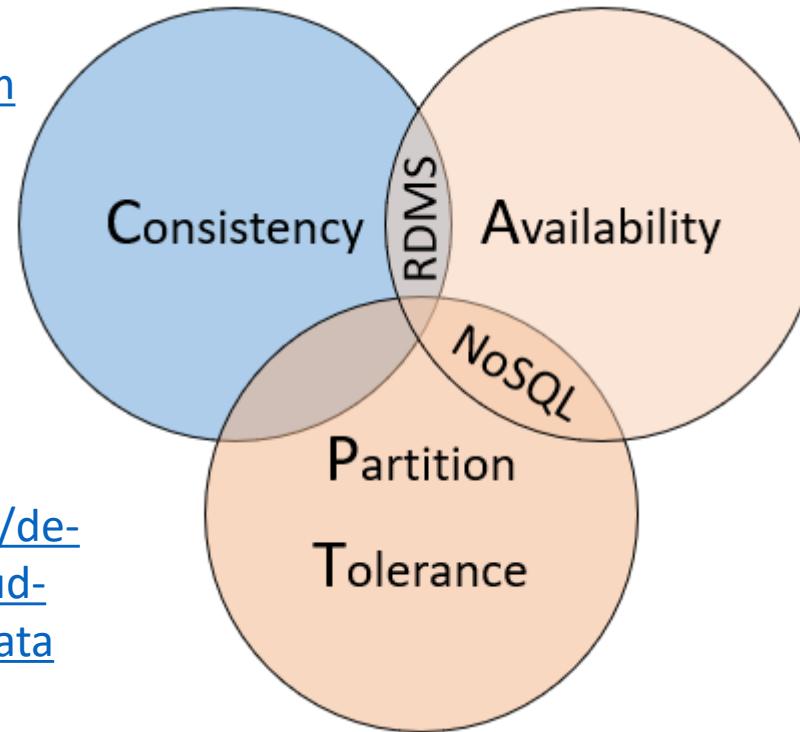
[Wie der Datentyp die Leistung beeinflussen kann.](#)

# CAP-Theorem

<https://www.ibm.com/de-de/cloud/learn/cap-theorem>

<https://tirsus.com/cap-theorem/>

<https://learn.microsoft.com/de-de/dotnet/architecture/cloud-native/relational-vs-nosql-data>



<https://datascience.eu/de/programmierung/cap-theorem-und-verteilte-datenbank-verwaltungssysteme/>

<http://www.webfunktion.de/cap-theorem/>

<https://blog.jetbrains.com/blog/2021/06/03/big-data-world-part-5-cap-theorem/>

<https://pol.techtec.world/de/blockchain/blockchain-cryptocurrency/captheorem>

## Die Theorie hinter dem ACID-Prinzip

Das **Theorie hinter dem ACID-Prinzip** setzt sich aus den folgenden Bausteinen zusammen:

- Das A (**atomacy**) in **ACID**: Man spricht dann von atomaren (**atomacy**) Operationen, wenn eine Sequenz von Datei-Operationen entweder ganz oder gar nicht ausgeführt wird.
- Das C (**consistency**) in **ACID**: Man spricht von einer vorhandenen Datenkonsistenz (**consistency**), wenn nach einer Sequenz von Datei-Operationen der Datenzustand in einem konsistenten Zustand hinterlassen wird.
- Das I (**isolation**) in **ACID**: Die Isolation (**isolation**) verhindert, dass sich parallele Ausführungen auf befindliche Datei-Operationen gegenseitig beeinflussen können.
- Das D (**durability**) in **ACID**: Die Dauerhaftigkeit (**durability**) gewährleistet, dass die Datei-Operationen dauerhaft auf einem Datenträger gesichert sind.

## Ein ACID-Prinzip Beispiel aus der Praxis

Um das **ACID-Prinzip** an einem **praxisnahen Beispiel zu erläutern**, werden wir den Verlauf einer Banküberweisung zwischen zwei Kreditinstituten beschreiben. Für diesen Vorgang sind zwei Vorgänge und drei Parteien nötig:

- Kreditinstitut A werden 100 Euro abgebucht.
- Kreditinstitut B werden 100 Euro gutgeschrieben.
- Kreditinstitut C

<https://databasecamp.de/daten/acid>

<https://www.bigdata-insider.de/was-ist-acid-a-776182/>

BASE ist ein Akronym und steht für

**"Basically Available, Soft-State, Eventually Consistent".**

Es beschreibt eine Klasse von Datenbank-Systemen, die sich von den traditionellen, sogenannten ACID-Systemen (Atomar, Konstant, Isoliert, Dauerhaft) unterscheiden.

<https://wi-wiki.de/doku.php?id=bigdata:konsistenz>

<https://wikis.gm.fh-koeln.de/Datenbanken/BASE>



# Rechtliche Grundlagen GDPR



<https://www.sailpoint.com/de/identity-library/was-ist-die-gdpr/>

<https://gdpr.eu/?cn-reloaded=1>

[https://www.haufe.de/compliance/management-praxis/datensicherheit/unterschied-zwischen-datenschutz-und-datensicherheit\\_230130\\_483954.html](https://www.haufe.de/compliance/management-praxis/datensicherheit/unterschied-zwischen-datenschutz-und-datensicherheit_230130_483954.html)

<https://www.brandmauer.de/blog/it-security/was-ist-der-unterschied-zwischen-datenschutz-und-datensicherheit>

<https://www.it-sicherheit-in-der-wirtschaft.de/ITS/Navigation/DE/Themen/Datenschutz-und-Datensicherheit/datenschutz-und-datensicherheit.html>

<https://www.jurando.de/datenschutz-datensicherheit-unterschied.php>

# DSGVO



<https://de.wikipedia.org/wiki/Datenschutz-Grundverordnung>

<https://dsgvo-gesetz.de/>

<https://dsiq.de/ds-gvo/>

# Bundesdatenschutzgesetz



<https://de.wikipedia.org/wiki/Bundesdatenschutzgesetz>

[https://www.iitr.de/blog/datenschutz-im-web-2-wann-gilt-das-deutsche-bundesdatenschutzgesetz/5518/?mtm\\_source=google&mtm\\_kwd=&mtm\\_medium=paid&mtm\\_campaign=perf-max-self-check&mtmclid=EAlaIQobChMI4oyTI5DR\\_AIVIKnVCh1TrgqVEAAYAiAAEgL-RfD\\_BwE&gclid=EAlaIQobChMI4oyTI5DR\\_AIVIKnVCh1TrgqVEAAYAiAAEgL-RfD\\_BwE](https://www.iitr.de/blog/datenschutz-im-web-2-wann-gilt-das-deutsche-bundesdatenschutzgesetz/5518/?mtm_source=google&mtm_kwd=&mtm_medium=paid&mtm_campaign=perf-max-self-check&mtmclid=EAlaIQobChMI4oyTI5DR_AIVIKnVCh1TrgqVEAAYAiAAEgL-RfD_BwE&gclid=EAlaIQobChMI4oyTI5DR_AIVIKnVCh1TrgqVEAAYAiAAEgL-RfD_BwE)

<https://dsgvo-gesetz.de/bdsg/>