

# Skalierbare NoSQL- und Cloud-Datenbanken in Forschung und Praxis

Felix Gessert, Norbert Ritter

Database and Information Systems Group  
Universität of Hamburg  
Vogt-Kölln Straße 33  
22527 Hamburg, Deutschland  
{gessert, ritter}@informatik.uni-hamburg.de

**Abstract:** Die rasante Entwicklung nicht-relationaler, verteilter NoSQL Datenbanksysteme hat in den letzten Jahren einen beispiellosen Aufschwung erlebt. Zwei zentrale Probleme haben diesen Prozess angestoßen: die gewaltigen Mengen von „User-generated-Content“ in modernen Anwendungen und die damit einhergehenden Anfragelasten und Datenvolumina, sowie der Ruf von Entwicklern nach problemspezifischen Datenmodellen und Schemaflexibilität. Das Tutorium „Skalierbare NoSQL- und Cloud-Datenbanken in Forschung und Praxis“ bietet einen umfassenden Überblick über die Konzepte und Techniken der relevantesten NoSQL- und Cloud-Datenbanken, mit einem besonderen Fokus auf Trade-Offs im Bereich Skalierbarkeit, Konsistenz und Anfragemächtigkeit. Es werden sowohl zugrundeliegende theoretische Konzepte, bahnbrechende wissenschaftliche Beiträge, als auch praktische Aspekte wie APIs, Lizenz- und Pricingmodelle sowie Architekturen diskutiert.

## 1 Beschreibung

Für einen breiten Überblick beginnt das Tutorium mit einer Motivation der NoSQL-Bewegung [SF12], sowie einer Kategorisierung von NoSQL-Systemen anhand ihres Datenmodells. Anschließend werden die wichtigen theoretischen Grundlagen eingeführt, wie beispielsweise Replikation, das CAP-Theorem [GL02] und Consistent Hashing. Diese Techniken sind die Grundlagen dafür, die Garantien und Fähigkeiten der anschließend detailliert beleuchteten Systeme beurteilen zu können. Zu den behandelten Systemen zählen u.a. die Dokumentendatenbank MongoDB, die Wide-Column Stores HBase und Cassandra, sowie die Key-Value Stores Redis und Riak. Anhand der Verteilungsarchitektur dieser Systeme mit ihren speziellen Replikations- und Partitionierungsmechanismen werden die erzielbaren Konsistenz- und Skalierbarkeitseigenschaften abgeleitet [FWGR14]. Um ein Verständnis für die jeweiligen Sweetspots im praktischen Einsatz zu erlangen, werden anschließend die Datenmodellierung und die Anfrageschnittstellen behandelt. Im Fokus stehen hier das Konzept der Schemalosigkeit und seine Implikationen, Key-Design, Schlüsselwahl bei Range- und Hash-Partitionierung, Query-Sprachen und Denormalisierungsstrategien.

Da Big Data Analytics und Data Science ein zunehmend wichtiges und stark diskutiertes Thema sind, werden auch Architekturen und Schnittstellen zur Integration operativer NoSQL-Datenbanken mit analytischen Systemen wie Hadoop und Spark diskutiert. Hierbei soll klar herausgearbeitet werden, dass Big Data zwei Ausprägungen hat: Data Management und Data Analytics, die im Kern jedoch auf denselben Abstraktionen aufbauen, um Skalierbarkeit zu erzielen und denselben Verfügbarkeits-Konsistenz-Trade-Offs unterworfen sind. Als Leitfaden und Orientierungshilfe werden die diskutierten NoSQL Datenbanken abschließend bezüglich mehrerer Eigenschaften wie Konsistenz, Anfragemächtigkeit und Lernkurve eingeordnet.

Der darauffolgende Abschnitt beschäftigt sich mit Cloud-Datenbanken [LS13] in seinen verschiedenen Ausprägungsformen: Cloud-hosted databases, Database-as-a-Service (DBaaS) [GBR14] und Backend-as-a-Service (BaaS) [GFW<sup>+</sup>14]. Einleitend wird die Multi-Tenancy-Problematik diskutiert, sowie das Konzept der Service-Levels-Agreements und des Workload-Managements eingeführt. Anhand repräsentativer Vertreter der jeweiligen Cloud-Datenbank-Kategorien (z.B. dem Amazon Relational Database Service als Beispiel für ein relationales DBaaS) soll hier ein praktischer Einblick in die konkrete Arbeit mit diesen Diensten gegeben werden. Da ein Großteil aller Cloud-Datenbank proprietärer Natur ist, müssen Rückschlüsse über nichtfunktionale Eigenschaften i.d.R. über vorhandene wissenschaftliche Publikationen (z.B. MegaStore [BBC<sup>+</sup>11] als Grundlage von Googles DataStore) angestellt werden. Hierbei werden wir zeigen, dass viele der Systeme dieselben Architekturen und Trade-Offs wie die zuvor diskutierten Open-Source NoSQL-Systeme aufweisen.

Als ein Beispiel dafür, wie NoSQL- und Cloud-Datenbanken auch althergebrachte Herangehensweisen in der Softwareentwicklung in Frage stellen, werden wird das Backend-as-a-Service Paradigma diskutieren, das es sich zum Ziel gemacht hat, die Web und Appentwicklung durch eine Kombination aus standardisierter Backendlogik und Cloud-Datenbanken radikal zu vereinfachen. Das Tutorium schließt mit einer Zusammenfassung und einer Aufstellung der wichtigsten offenen Forschungsfragen für NoSQL- und Cloud-Datenbanken.

## 2 Zielgruppe

Das Tutorium richtet sich an Anwender (z.B. Softwareentwickler und -Architekten) und Wissenschaftler, die ihre Kenntnisse im Bereich NoSQL- und Cloud-Datenbanken vertiefen oder sich über die wichtigsten offenen Probleme informieren möchten. Vorkenntnisse in den Bereichen Datenbanksysteme, verteilte Systeme und Softwareentwicklung sind von Vorteil, aber nicht zwingend erforderlich.

### 3 Organisatoren

**Felix Gessert** ist Doktorand an der Gruppe Datenbanken und Informationssysteme der Universität Hamburg. Er forscht in den Bereichen skalierbare Datenbanksysteme, Transaktionen, Cloud Computing und Web-Technologien. Sein besonderer Fokus gilt dem Caching und der Transaktionsverarbeitung im Kontext von Cloud-Datenbanksystemen. Er ist außerdem Gründer des Spin-Offs Baqend, das die Forschungsergebnisse in Form einer Backend-as-a-Service Plattform umsetzt.

**Norbert Ritter** ist Professor für Informatik und Leiter der Gruppe Datenbanken und Informationssysteme. Seine Forschungsinteressen umfassen u.a. verteilte und föderierte Datenbanksysteme, Transaktionsverarbeitung, Caching, Cloud Data Management, Informationsintegration und autonome Datenbanksysteme.

### Literatur

- [BBC<sup>+</sup>11] Jason Baker, Chris Bond, James C. Corbett, J. J. Furman, Andrey Khorlin, James Larson, Jean-Michel Léon, Yawei Li, Alexander Lloyd und Vadim Yushprakh. Megastore: Providing Scalable, Highly Available Storage for Interactive Services. In *CIDR*, Jgg. 11, Seiten 223–234, 2011.
- [FWGR14] Steffen Friedrich, Wolfram Wingerath, Felix Gessert und Norbert Ritter. NoSQL OLTP Benchmarking: A Survey. In Erhard Plödereder, Lars Grunske, Eric Schneider und Dominik Ull, Hrsg., *44. Jahrestagung der Gesellschaft für Informatik, Informatik 2014, Big Data - Komplexität meistern, 22.-26. September 2014 in Stuttgart, Deutschland*, Jgg. 232 of *LNI*, Seiten 693–704. GI, 2014.
- [GBR14] Felix Gessert, Florian Bucklers und Norbert Ritter. Orestes: A scalable Database-as-a-Service architecture for low latency. In *Workshops Proceedings of the 30th International Conference on Data Engineering Workshops, ICDE 2014, Chicago, IL, USA, March 31 - April 4, 2014*, Seiten 215–222. IEEE, 2014.
- [GFW<sup>+</sup>14] Felix Gessert, Steffen Friedrich, Wolfram Wingerath, Michael Schaarschmidt und Norbert Ritter. Towards a Scalable and Unified REST API for Cloud Data Stores. In Erhard Plödereder, Lars Grunske, Eric Schneider und Dominik Ull, Hrsg., *44. Jahrestagung der Gesellschaft für Informatik, Informatik 2014, Big Data - Komplexität meistern, 22.-26. September 2014 in Stuttgart, Deutschland*, Jgg. 232 of *LNI*, Seiten 723–734. GI, 2014.
- [GL02] Seth Gilbert und Nancy A. Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News*, 33(2):51–59, 2002.
- [LS13] Wolfgang Lehner und Kai-Uwe Sattler. *Web-Scale Data Management for the Cloud*. Springer, 2013.
- [SF12] Pramod J. Sadalage und Martin Fowler. *NoSQL distilled: a brief guide to the emerging world of polyglot persistence*. Pearson Education, 2012.

