

# 1 Einführung

## 1.1 Motivation

In den letzten Jahren ist für viele Anwendungen das Suchfeld immer wichtiger geworden. War die Suche früher noch ein Teilaspekt einer Anwendung, wird heutzutage immer mehr Funktionalität über Suchserver wie Elasticsearch abgewickelt. Viele Benutzer greifen auf Web-Inhalte nur noch über die Suche zu, das Menü als Hilfsmittel für die Navigation verliert an Bedeutung. Dabei soll alles möglichst so gut auffindbar sein, wie man es von Diensten wie Google gewohnt ist.

Gleichzeitig können Suchserver auch zur Lösung anderer Probleme genutzt werden. Sie sind oftmals für lesenden Zugriff optimiert und können skalieren. So bieten sie sich für Aufgaben an, die klassischerweise von einer Datenbank übernommen wurden. Anwendungen verkraften dadurch hohe Zugriffszahlen, Lastspitzen können abgefangen werden.

Durch die Verfügbarkeit von flexiblen Suchservern wie Elasticsearch ergeben sich jedoch auch ganz neue Anwendungsfälle – das zentralisierte Logging, bei dem alle anfallenden Lognachrichten einer Anwendung in einem zentralen Speicher zur Auswertung geschrieben werden, kann Zeit sparen und Entwicklernerven schonen.

Neben den Anwendungs-Logs kann das Ganze dann auch noch auf Business-Daten übertragen werden. Durch die flexiblen Auswertungsmöglichkeiten können Daten analysiert und visualisiert werden. Dabei kann der Suchserver bestehende Data-Warehouse-Lösungen ergänzen oder gar ersetzen.

Bevor wir uns in einem ersten Praxisbeispiel anschauen, wie sich die Nutzung von Elasticsearch anfühlt, werfen wir noch einen Blick auf die Geschichte von Elasticsearch, die stark mit der zugrunde liegenden Bibliothek Apache Lucene zusammenhängt.

## 1.2 Geschichte von Elasticsearch

Wenn es um die Volltextsuche geht, kam man lange Jahre nicht an der Java-Bibliothek Lucene vorbei, die 1999 von Doug Cutting entworfen und bei SourceForge unter einer Open-Source-Lizenz veröffentlicht wurde. Seit 2001 wird das

Projekt bei der Apache Software Foundation weiterentwickelt. Lucene implementiert einen sogenannten invertierten Index, der die Basis für viele Suchlösungen bildet. Damit können Daten performant und flexibel durchsucht werden. Neben der Implementierung auf der JVM in Java entstanden über die Jahre auch noch unterschiedliche Portierungen in anderen Programmiersprachen, beispielsweise PyLucene für Python oder Lucene.NET für die .NET-Runtime [1].

Einen Wandel in der Nutzung gab es in den Jahren nach 2006, als der hauptsächlich von Yonick Seeley geschriebene Suchserver Solr von CNET als Open-Source-Projekt an die Apache Software Foundation übergeben wurde. Apache Solr baut auf Lucene auf und stellt die zugrunde liegende Funktionalität über eine HTTP-Schnittstelle zur Verfügung. [2] Damit kann auch außerhalb der JVM die Funktionalität von Lucene genutzt werden. Durch weitere Features wie die ursprünglich nur in Solr verfügbare Facettierung und den konfigurativen Ansatz wurde Solr schnell ein großer Erfolg.

Im Jahre 2010 schließlich veröffentlichte Shay Banon die erste Version von Elasticsearch als Open-Source-Projekt unter der Apache-2.0-Lizenz. Viele der Funktionen basieren auf seinen Erfahrungen bei der Entwicklung des Lucene-Such-Frameworks Compass [3]. Auf den ersten Blick sieht Elasticsearch dabei ähnlich wie Apache Solr aus – es baut ebenfalls auf Lucene auf und stellt die Funktionalität per HTTP zur Verfügung. Im Detail gibt es jedoch große Unterschiede, etwa eine kompromisslose Fokussierung auf Verteilung. Durch die transparente Datenverteilung und die einfachen, modernen Bedienkonzepte wurde Elasticsearch schnell erfolgreich und ist mittlerweile nicht nur bei jüngeren Unternehmen wie GitHub oder Stackoverflow im Einsatz, sondern auch in konservativeren wie Banken oder unterschiedlichen amerikanischen Behörden. Seit 2012 existiert mit elastic (vormals Elasticsearch) ein Unternehmen, das Support und kommerzielle Erweiterungen für Elasticsearch anbietet.

Elasticsearch stellt viele unterschiedliche Funktionen zur Verfügung. Dazu gehören unterschiedliche Möglichkeiten zur Indizierung von Inhalten und der Suche darauf sowie unterstützende Funktionalitäten wie die Relevanzberechnung oder Möglichkeiten zur automatischen Vervollständigung von Inhalten. Einzelne Elasticsearch-Instanzen bilden einen Cluster und Daten können in einzelne Shards aufgeteilt werden, die auf unterschiedliche Knoten verteilt werden können. Durch die verteilte Natur ist es deutlich einfacher, als mit vielen anderen Systemen auch mit großen Datenmengen umzugehen, was beispielsweise für Anwendungsfälle wie das zentralisierte Logging wichtig ist. Aggregationen können schließlich neue Blickwinkel auf Daten ermöglichen.

Dieses Buch soll einen Einblick in die Nutzung von Elasticsearch geben und zeigen, wie der Suchserver für unterschiedliche Anwendungen zum Einsatz kommen kann. Der Start mit Elasticsearch fällt glücklicherweise ganz leicht, wie wir im folgenden Beispiel sehen, bei dem wir Elasticsearch installieren, Daten darin speichern und direkt durchsuchen.

## 1.3 Ein erstes Beispiel

Voraussetzung für die Installation von Elasticsearch ist lediglich eine aktuelle Java VM, die unter <http://java.com> für unterschiedliche Systeme heruntergeladen werden kann. Auf der Elasticsearch-Homepage stehen unter <https://www.elastic.co/downloads/elasticsearch> unterschiedliche Archive zur Verfügung, die direkt entpackt werden können. In diesem Buch wird durchgehend Elasticsearch in der Version 1.6.0 eingesetzt. Elasticsearch kann nach dem Entpacken über ein Skript gestartet werden<sup>1</sup>.

```
wget https://download.elastic.co/elasticsearch
↪/elasticsearch/elasticsearch-1.6.0.zip2
unzip elasticsearch-1.6.0.zip
elasticsearch-1.6.0/bin/elasticsearch
```

Die Funktionalität steht über eine HTTP-Schnittstelle zur Verfügung. Ob die Anwendung gestartet wurde, sehen wir nicht nur an den Logausgaben, sondern auch durch einen Zugriff auf den Server per HTTP, etwa über das Kommandozeilenwerkzeug `cURL`.

```
curl -XGET "http://localhost:9200"
```

### Elasticsearch und HTTP

Da HTTP mittlerweile allgegenwärtig ist, können unterschiedliche Werkzeuge für den Zugriff verwendet werden. In diesem Buch wird wie in vielen anderen Dokumentationsquellen der Kommandozeilen-Client `cURL` [4] verwendet, der für unterschiedliche Systeme zur Verfügung steht. Es ist jedoch auch problemlos eine Nutzung anderer Werkzeuge möglich, beispielsweise von Browser-Plugins wie Postman [5] für Chrome oder `RESTClient` [6] für Firefox.

Die Kommunikation mit Elasticsearch über HTTP wird oft als RESTful bezeichnet. REST ist ein Architekturstil, bei dem die Funktionalität einer Anwendung über Ressourcen abgebildet wird, die über HTTP angesprochen werden können. Ressourcen sind durch URIs identifizierbar und können unterschiedliche Repräsentationen haben. Zustände werden über Links zwischen den Ressourcen modelliert. [7]

Die Diskussion, ob die Schnittstellen von Elasticsearch als RESTful anzusehen sind, wird in diesem Buch nicht geführt. Wichtig ist, dass die Nutzung von HTTP zur Kommunikation einen einfachen Umgang mit Elasticsearch sowohl für den Betrieb als auch für die Entwicklung erlaubt.

---

<sup>1</sup>Im Buch wird Elasticsearch in einer Linux-Umgebung verwendet. Ein Betrieb ist jedoch auch problemlos unter Windows möglich.

<sup>2</sup>Lange Ein- und Ausgaben auf der Kommandozeile, die eigentlich durchgängig zu schreiben sind, werden aus Platzgründen im Buch gelegentlich auf mehrere Zeilen umbrochen. Dies ist dann mit dem Zeichen `↪` gekennzeichnet.

### Sense

Gerade für Einsteiger kann neben der Nutzung eines einfachen Tools zur Kommunikation über HTTP auch der Einsatz von Marvel [8] eine Alternative sein, einem Monitoring-Plugin für Elasticsearch. Dieses integriert das Werkzeug Sense, mit dem komfortabel Abfragen formuliert werden können. Neben der kontextsensitiven Autovervollständigung von Query-Bestandteilen bietet das Werkzeug Syntax-Highlighting, eine History vergangener Abfragen und die Möglichkeit, Abfragen von und in cURL-Aufrufe umzuwandeln.

Bei Marvel handelt es sich um ein kostenpflichtiges Plugin, das allerdings zur Entwicklungszeit kostenlos genutzt werden kann. Auf [elasticsearch-buch.de](http://elasticsearch-buch.de) findet sich eine kurze Einführung zur Nutzung von Sense in Marvel.

Wenn die Anwendung erfolgreich gestartet wurde, werden einige Informationen angezeigt.

```
{
  "status" : 200,
  "name" : "Wolf",
  "cluster_name" : "elasticsearch",
  "version" : {
    "number" : "1.6.0",
    "build_hash" : "cdd3ac4dde4f69524ec0a14de3828cb95bbb86d0",
    "build_timestamp" : "2015-06-09T13:36:34Z",
    "build_snapshot" : false,
    "lucene_version" : "4.10.4"
  },
  "tagline" : "You Know, for Search"
}
```

Elasticsearch verwendet zur ein- und ausgehenden Kommunikation ausschließlich JSON-Dokumente, im Beispiel ein Dokument, das Informationen zum Server wie den Namen, die Version und Informationen zur verwendeten Lucene-Version zusammenfasst. Zur Speicherung der Daten können ebenfalls JSON-Dokumente verwendet werden, die beliebig aufgebaut sein können. Über eine POST-Anfrage kann ein neues Dokument hinzugefügt werden.

```
curl -XPOST "http://localhost:9200/example/doc" -d'
{
  "title": "Hallo Welt",
  "tags": ["example", "elasticsearch"]
}'
```

Diese gespeicherten Daten können im Anschluss sofort durchsucht werden.

```
curl -XGET "http://localhost:9200/_search?q=welt&pretty"
```

Als Ergebnis erhalten wir für den Suchbegriff *welt* unser gerade hinzugefügtes Dokument zurück.

```
{
  "took" : 231,
  [...],
  "hits" : {
    "total" : 1,
    "max_score" : 0.15342641,
    "hits" : [ {
      "_index" : "example",
      "_type" : "doc",
      "_id" : "AU3g6jivOZtCV272fVut",
      "_score" : 0.15342641,
      "_source":{
        "title": "Hallo Welt",
        "tags": ["example", "elasticsearch"]
      }
    } ]
  }
}
```

Wir haben in wenigen Schritten eine Möglichkeit geschaffen, unsere Daten durchsuchen zu können.

#### **Formatierung der Ausgabe**

Viele der JSON-Dokumente, die von Elasticsearch zurückgegeben werden, sind nicht formatiert und deshalb schwer lesbar. Über den optionalen Parameter `pretty`, der an die URL angehängt werden kann, wird die Rückgabe formatiert. Im Buch werden die Ausgaben oftmals formatiert dargestellt, der Parameter wird aus Gründen der Lesbarkeit aber nicht explizit angegeben.

## 1.4 Anwendungsfälle

Elasticsearch kann für die unterschiedlichsten Einsatzgebiete nützlich sein. Oftmals wenden Anwendungen unterschiedliche Aspekte an und kombinieren diese. Zu den wichtigsten Einsatzszenarien gehören die folgenden.

**Volltextsuche** Durch den Unterbau Apache Lucene ist Elasticsearch prädestiniert zur Volltextsuche. Dabei kann innerhalb von Texten nach einzelnen Schlüsselwörtern, Phrasen oder unvollständigen Begriffen gesucht werden. Gerade auf Webseiten oder bei Anwendungen, die große Datenmengen zur Verfügung stellen, kann die Volltextsuche ein elementares Werkzeug für die Nutzer sein.

**Abfrage strukturierter Daten** Neben der Suche in Textinhalten können in Elasticsearch Daten auch strukturiert abgefragt werden. Es werden unterschiedliche Datentypen für numerische Daten, Datumswerte und sogar Geodaten unterstützt, die auf unterschiedliche Art, auch kombiniert, abgefragt werden können.

**Analytics** Mit den Aggregationen hat Elasticsearch ein mächtiges Werkzeug für Analytics integriert. Damit können Besonderheiten und Gemeinsamkeiten in Datensätzen auch aus großen Datenmengen performant und trotzdem flexibel extrahiert werden. Durch die Sharding-Funktionalität können auch sehr große Datenmengen in Echtzeit gespeichert, durchsucht und aggregiert werden.

**Verarbeitung hoher Leselast** Elasticsearch kann durch die horizontale Skalierung gut mit vielen Leseanfragen umgehen. Daten können auf unterschiedliche Knoten repliziert und damit auch eine hohe Abfragelast verarbeitet werden.

**Datenkonsolidierung** Durch die einfache Dokumentenstruktur können Daten aus unterschiedlichen Quellen in Elasticsearch zusammengeführt und gemeinsam abgefragt werden.

**Logfile-Analyse** Durch die Verteilung unserer Anwendungen in unterschiedliche Prozesse und auf unterschiedliche Maschinen ist die Analyse im Fehlerfall oder auch die Kontrolle der Logdaten deutlich komplexer geworden. Elasticsearch ist vor allem in Verbindung mit den Werkzeugen Logstash und Kibana sehr beliebt als zentraler Datenspeicher für Log-Events.

## 1.5 Wann Elasticsearch?

Neben den unterschiedlichen Anwendungsfällen bietet Elasticsearch noch weitere Vorteile.

**Einfacher Start** Die Nutzung von Elasticsearch gestaltet sich auch für Einsteiger sehr einfach. Es ist wenig Konfiguration notwendig, Elasticsearch arbeitet in vielen Fällen mit sinnvollen Standardwerten.

**Programmiersprachunabhängigkeit** Durch die Nutzung von HTTP und JSON kann Elasticsearch aus so gut wie jeder Programmiersprache angesprochen werden. Zusätzlich existieren zahlreiche Client-Bibliotheken für die verbreitetsten Programmiersprachen.

**Open-Source-Projekt** Da Elasticsearch unter der bewährten Apache-Lizenz veröffentlicht wird, ist der Einsatz im kommerziellen Umfeld kein Problem. Zusätzlich kann der Quelltext genutzt werden, um die Funktionsweise zu verstehen und auch um eigene Fehlerbehebung oder Erweiterungen durchzuführen.

Eine Nutzung von Elasticsearch sollte in vielen Fällen keine Probleme darstellen. Jedoch sollen auch noch die folgenden Aspekte erwähnt werden, die unter Umständen für andere Systeme sprechen können.

**Stärke durch Flexibilität** Elasticsearch spielt seine Stärke bei den flexiblen Abfragen aus. Wenn immer nur über eine ID auf die Dokumente zugegriffen werden soll, können andere Systeme wie Apache Cassandra eine bessere Wahl sein.

**Komplexität** Elasticsearch bringt als von Grund auf verteiltes System inhärente Komplexität mit sich. Vieles davon ist vor dem Anwender versteckt, im Ernstfall muss man sich jedoch mit unterschiedlichen Eigenheiten auseinandersetzen. Wenn die Verteilung nicht benötigt wird, kann eine Nutzung anderer Systeme einfacher sein.

**Projekt wird von einer Firma getrieben** Im Gegensatz zu Projekten beispielsweise bei der Apache Software Foundation steht hinter Elasticsearch eine Firma, die die Entwicklung steuert. Dies muss kein Nachteil sein, das Unternehmen elastic stellt viele exzellente Entwickler zur Arbeit an Elasticsearch und auch an Apache Lucene ab. Potenziell hat man jedoch mit einem Projekt bei einer Stiftung mehr Sicherheit, was eine mögliche Einflussnahme in die Entwicklung betrifft.

**Weniger gut geeignet für Daten mit hoher Änderungshäufigkeit** Elasticsearch und das darunter liegende Lucene gewinnen viel der Geschwindigkeit dadurch, dass einmal geschriebene Daten möglichst nicht mehr verändert werden. Wenn durch die Anwendung bestimmt wird, dass sich bestehende Datensätze sehr häufig ändern sollen, können andere Systeme zur Speicherung besser geeignet sein.

Auch wenn Elasticsearch für viele Anwendungen empfohlen werden kann, ist je nach den eigenen Anforderungen eine Evaluierung unterschiedlicher Systeme sinnvoll.

## 1.6 Über dieses Buch

Dieses Buch soll den neugierigen Leser in die Grundlagen von Elasticsearch einführen. Dabei werden unterschiedliche Einsatzzwecke, vor allem die Volltextsuche, allerdings auch die Nutzung für Analytics oder als Datenspeicher für zentralisiertes Logging beleuchtet.

Zum Verständnis der Konzepte und Beispiele sind grundlegende Kenntnisse in der Softwareentwicklung und erste Erfahrungen in der Datenmodellierung, beispielsweise für relationale Datenbanksysteme, sinnvoll. Das Buch ist so gestaltet, dass die Kapitel in der angegebenen Reihenfolge gelesen werden. Für Leser, die sich bereits mit Elasticsearch auskennen, können auch einzelne Kapitel als Referenz dienen.

**Kapitel 2** zeigt, wie klassischerweise eine Suchanwendung entsteht. Dabei werden unterschiedliche Konzepte wie Indizierung, Mapping, Facettierung über Aggregationen und die Query-DSL erläutert.

**Kapitel 3** widmet sich dem Umgang mit Texten. Dabei wird auf unterschiedliche Mechanismen eingegangen, die viel zum Eindruck einer intelligenten Suche beitragen. Einerseits werden der Umgang mit unterschiedlichen Sprachen als auch andererseits nützliche Funktionen wie die Hervorhebung vom Suchbegriff oder Autovervollständigung erläutert.

**Kapitel 4** hat die Relevanzberechnung zum Thema, ein großes Unterscheidungsmerkmal zu datenbankbasierten Systemen. Gerade bei großen Datenmengen kann ein Nutzer oftmals sehr viele Ergebnisse erhalten – die Relevanzsortierung hilft dann, die besten Treffer direkt nach vorne zu sortieren. Über unterschiedliche Mechanismen wie Boosting oder die Function-Score-Query kann die Relevanzsortierung auch für eigene Anwendungen beeinflusst werden.

**Kapitel 5** In Kapitel 5 geht es um die Indizierung. Dabei werden einige Strategien und Mechanismen zur Anbindung von externen Datenquellen erläutert. Um zu verstehen, wie und wann Dokumente persistiert werden, ist ein Exkurs in einige Lucene-Interna notwendig.

**Kapitel 6** dreht sich um alle Aspekte, die für die Verteilung der Daten wichtig sind. Im Zentrum stehen Shards und Replicas und wie diese auf die unterschiedlichen Knoten im Cluster verteilt werden.

**Kapitel 7** erläutert unterschiedliche Mechanismen, wie Daten für Elasticsearch modelliert werden. Dazu gehören die elementaren Datentypen, die Gestaltung der Indexstruktur und die Verwaltung von Beziehungen zwischen Inhalten.

**Kapitel 8** stellt die unterschiedlichen Aggregationen vor, die genutzt werden können, um Daten anhand bestimmter Eigenschaften zusammenzufassen und um Berechnungen darauf durchzuführen.

**Kapitel 9** beschreibt den Zugriff auf Elasticsearch und die Nutzung des Java- und JavaScript-Clients.

**Kapitel 10** hat einige für den Betrieb wichtige Themen zum Inhalt. Dazu gehören einerseits Überlegungen, die vor einer Produktivnahme anstehen, als auch Themen wie Monitoring oder Datensicherung.

**Kapitel 11** widmet sich schließlich der populären Nutzung von Elasticsearch zur Speicherung von Logdaten mittels der Werkzeuge Logstash und Kibana, die zusammen den ELK-Stack ausmachen. Zusätzlich wird die Alternative Graylog vorgestellt, eine integrierte Lösung, die ebenfalls auf Elasticsearch zur Datenspeicherung setzt.

**Kapitel 12** gibt einen Ausblick in die Zukunft von Elasticsearch und stellt in Kürze einige Themen vor, auf die im Buch nicht eingegangen werden konnte.

In den Anhängen finden sich schließlich noch Informationen zu einem Vorgehen zum Neuindizieren von Dokumenten und Hinweise zur Installation des Twitter-Rivers.

Die Beispiele und Erläuterungen in diesem Buch bauen auf der Elasticsearch-Version 1.6.0 auf. Aus Platzgründen musste die Darstellung von Ein- und Ausgaben teilweise umbrochen und abgekürzt werden. Dadurch sind die Codestücke



eventuell nicht immer so wie sie sind, lauffähig. Alle Beispiele in ihrer Originalform, weitere Informationen und mögliche Korrekturen finden sich auf der Homepage [elasticsearch-buch.de](http://elasticsearch-buch.de).

In den Kapiteln 2–5 sind einzelne Aufgaben verstreut, die mit »Zum Selbermachen« markiert sind und direkt zum praktischen Mitarbeiten anregen sollen. Hinweise und mögliche Lösungen zu diesen Aufgaben finden sich ebenfalls auf [elasticsearch-buch.de](http://elasticsearch-buch.de).

## 1.7 Danksagung

Dieses Buch konnte nur entstehen, weil ich auf die Hilfe vieler Leute zurückgreifen konnte. An erster Stelle möchte ich meinen Reviewern danken, die wertvollen Input zum Manuskript in unterschiedlichen Stadien gegeben haben. Vielen Dank an Tobias Kraft, Andreas Jägle, Bernd Fondermann, Eberhard Wolff, Jochen Schalanda und Dr. Patrick Peschlow.

Vielen Dank auch an das Team beim [dpunkt.verlag](http://dpunkt.verlag) um meinen Lektor René Schönfeldt und an Niko Köbler, der dafür gesorgt hat, dass ich dieses Buch schreiben konnte.

Ebenfalls vielen Dank an all die Leute, die ihr Wissen in Blogposts, auf Konferenzvorträgen oder im direkten Austausch weitergeben und damit indirekt mitgeholfen haben, dass dieses Buch entstehen konnte. Ohne Anspruch auf Vollständigkeit gehören dazu Alexander Reelsen (der mich auch ursprünglich von Elasticsearch überzeugt hat), Britta Weber, Adrien Grand, Uwe Schindler, Clinton Gormley, Christian Uhl, Jörg Prante, Andrew Cholakian, Michael McCandless, Zachary Tong, Lucian Precup und Alex Brasetvik.

Vielen Dank auch an alle Entwickler hinter Elasticsearch und Apache Lucene, die es ermöglichen, dass die Software für uns alle verfügbar ist.

Zum Schluss noch vielen Dank an meine Eltern, die mich jahrelang unterstützt haben, und vor allem an Lianna, die während der Entstehung dieses Buches oft auf mich verzichten musste, mir aber immer eine Hilfe war.