

NoSQL-tietokannat tieteellisen tutkimusaineiston arkistoinnissa

Marko Peltola

Jyväskylän yliopisto

17.3.2011

- Kandidaatintutkielma NoSQL-tietokannoista
⇒ teoriatausta
- Gradussa pääosassa tutkimusaineistojen arkistointi
⇒ tapaustutkimus / konstrukttiivinen tutkimus

Mikä NoSQL?

- NoSQL \neq !SQL
- Not Only SQL
- Ei-relaatiotietokanta

- Relaatitietokannat
(MySQL, PostgreSQL)
- Oliotietokannat
- Dokumenttitietokannat
(CouchDB, MongoDB)
- Avain–arvo-tietokannat
(Dynamo, Riak, Redis)
- Graafitietokannat
(Neo4j)
- Sarakeorientoituneet tietokannat
(Bigtable, Cassandra, HBase)

Miksi NoSQL?

- Big Data
- Pilvilaskenta
- Skaalautuvuus
- Ei kiinteää skeemaa
- Google, Facebook, Amazon, ...

Perinteiset tietokantajärjestelmät: ACID

ACID takaa, että tietokantatransaktiot suoritetaan luotettavasti

Atomicity (jakamattomuus)

Tapahtuma suoritetaan joko kokonaan tai ei lainkaan

Consistency (eheys)

Tapahtumien myötä tietokanta siirtyy ehestä tilasta eheään tilaan

Isolation (eristyneisyys)

Tapahtuman ollessa kesken ei muut tapahtumat saa käyttää sen ominaisuuksia

Durability (pysyvyys)

Vahvistetun tapahtuman on säilyttävä

Hajautetut järjestelmät: CAP-teoreema

Consistency (eheys)

Jokainen solmu näkee samat tiedot samaan aikaan

Availability (saatavuus)

Aina on pystyttävä kirjoittamaan ja lukemaan

Partition-tolerance (osituksen sietokyky)

Järjestelmä toimii, vaikka verkkoyhteys joidenkin solmujen välillä katkeaa

CAP-teoreeman mukaan hajautettu järjestelmä voi taata näistä kaksi, muttei kaikkia kolmea

Partition-tolerance

P

Bigtable
HBase
MongoDB
Redis

Cassandra
CouchDB
Dynamo
Riak

Valitse
kaksi

C

Consistency

A

Availability

MySQL
PostgreSQL

- Sovellusprojekti kevät 2011 (Judo)
- Laajojen tutkimusaineistojen arkistointi
- Kyettävä tallentamaan isoja tiedostoja (> 25 GB)
- Metadatat heterogeenisiä
- WWW-palvelu
- MongoDB + GridFS

- Kuinka hyvin valittu tietokanta soveltuu tutkimusaineistojen arkistointiin?
- Onko tiedostoja järkevä tallentaa tietokantaan? (GridFS)
- Metadatan muoto?

- Chang et al., "Bigtable: A Distributed Storage System for Structured Data", 2008
- DeCandia et al., "Dynamo: Amazon's Highly Available Key-value Store", 2007
- Lakshman, Malik, "Cassandra – A Decentralized Structured Storage System", 2010
- Leavitt, "Will NoSQL Databases Live Up to Their Promise?", 2010
- Stonebraker et al., "The end of an architectural era (it's time for a complete rewrite)", 2007