# Which databases solve my problem?

## a survey of open source databases

Selena Deckelmann
End Point Corporation
@selenamarie
PostgreSQL Global Development Group

LCA 2010
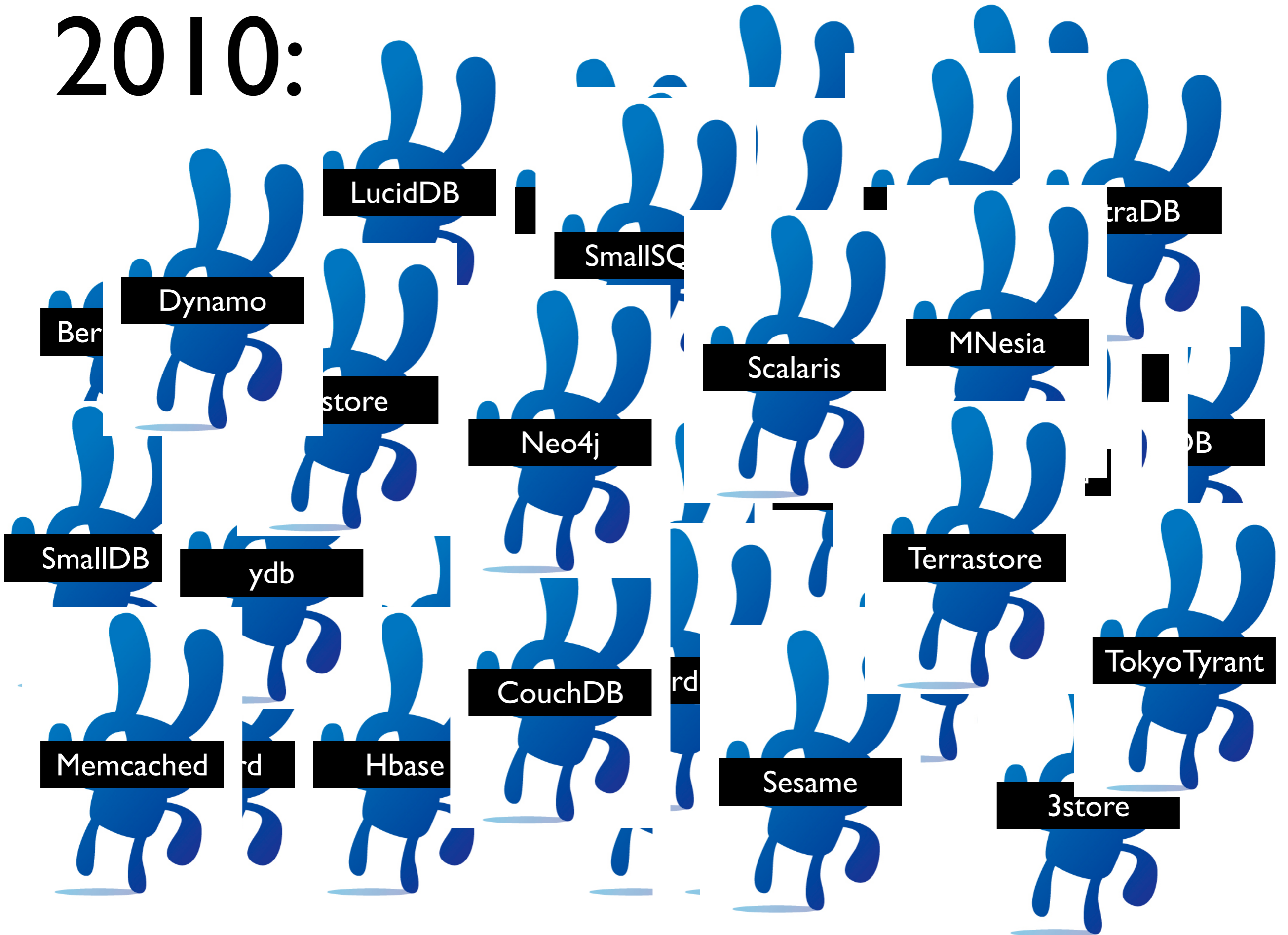
**END** *point*

# 2005:

BerkeleyDB

MySQL

PostgreSQL

SQLite

ENDpoint

# 2010:

LucidDB

traDB

SmallSQ

Dynamo

Ber

Scalaris

MNesia

store

Neo4j

DB

SmallDB

Terrastore

ydb

TokyoTyrant

CouchDB

rd

Memcached    rd    Hbase    Sesame    3store

2am on Monday morning.

# Which open source database should I use?

**END**_point_

via http://www.oddee.com/item_86516.aspx

# MySQL vs PostgreSQL

END*point*

# What problem are you trying to solve?

END*point*

# Some problems:

END*point*

# I need to store and manipulate GIS data.

END*point*

# I need a database for my blog.

I have ONE BILLION users to store and analyze data from.

END*point*

# Define your problem.



LCA 2010

END point

# Which problems are important?

END*point*

# performance

your use case.
test with real data.

END *point*

# interoperability
## can I get my data in/out?
## how painful is it?

END*point*

# sustainability
## how is the software made?

END*point*

# Which databases solve my problem?

**END**_point_

# The Survey

- Wasn't perfect.

- Contacted 25 projects, 12 responses.

- Will try again with different questions, cooler website.

END*point*

# The questions:

- What is the name of your project?
- How would you describe your software and what it does in a sentence or two?
- Who is the target user or audience for your database? Do you have any case studies to share?
- Is there a proprietary work-alike or equivalent to your open source database?
- What's the best mailing list for users of your database to subscribe to?
- What's the best mailing list for developers of your database to subscribe to?
- What's the best document for new developers to read if they want to get involved?

END*point*

- What revision control system does your project primarily use?
- What motivated you to create a new project, rather than join an existing project?
- Do you have a roadmap for the next year? If so, what is it?
- Does anyone provide commercial support for your software?
- What languages are drivers available in, and/or what protocols does your database support? Are they up to date?
- Do you need help with any particular drivers?
- Is there some question I should have asked?
- What feature(s) sets your project apart from your peers?

And I did my own research...

END *point*

# Means of comparison

Database model

Infrastructure features

Development style

END *point*

# Models: defining what operations you'll likely perform on the data

END*point*

# Relational Database models

**OLTP**: Transaction-oriented

**Embedded**: Bundling, simplicity, testing

**Column**: Data warehouses

**MPP**: Massively Parallel

**Streaming**: Query streams, not storage

END*point*

# Relational Database Models

| OLTP | Embedded | Column-store |
|---|---|---|
| CUBRID<br>MySQL (InnoDB)<br>PostgreSQL | H2<br>HSQLDB<br>SQLite | MonetDB<br>LucidDB<br>C-store/Vertica<br><br>(Cassandra<br>Hbase) |

END *point*

# non-Relational Database models

**Flatfile**: See Tin ( http://tr.im/KNFp )

**Key-value**: map-reduce, fault-tolerance, caching

**Multi-value**: Multi-dimensional - GT.M

**Graph/Triple-store**: Relationship queries

**Document-oriented**: Semi-structured data

END *point*

# non-Relational Database Models

| Key-value | Graph/ Triple-store | Document |
|---|---|---|
| BerkeleyDB Cassandra Hbase Memcached Riak Redis TokyoCabinet ydb | Neo4j 4store Parliament | CouchDB BerkeleyDB-XML MongoDB |

# infrastructure features:

"distributed"

memory

HA

# "Distributed"

| Partitioning/ Sharding | Replication | |
|---|---|---|
| Cassandra Hbase Voldemort Riak MySQL | BerkeleyDB CouchDB Cassandra MySQL PostgreSQL Riak | Scalaris Voldemort HyperTable HBase Memcached MNesia |

END *point*

# Memory vs Disk

| In-memory* | Configurable | Disk |
|---|---|---|
| Memcached<br>Scalaris<br>Redis | Cassandra<br>Hbase<br>HyperTable<br>MNesia | Everyone else |

*This is databases existing solely in memory and being unable or never persisting to disk.

LCA 2010

END *point*

# High Availability

| Node failover |
|---|
| Cassandra<br>HBase<br>Riak |

Otherwise, use one or more of: heartbeat, DBRD, filesystem replication, etc.

ENDpoint

# Code Development Model

| Core + modules | Monolithic | Infrastructure |
|---|---|---|
| Drizzle LucidDB PostgreSQL | GT.M Ingres CUBRID | Memcached Redis Scalaris |

LCA 2010

END *point*

# Community Development Model

| Benevolent Dictator | Feature driven | Small Group | A mix |
|---|---|---|---|
| Redis XtraDB MckoiDDB | Apache Derby InfiniDB SmallSQL | CouchDB MonetDB Riak | LucidDB Drizzle H2 PostgreSQL |

ENDpoint

# Plans for the data

- Attempt to update Wikipedia

- Talk to people who write real surveys

- Contacting more projects

- http://ossdbsurvey.org

END*point*

# The Future!

# Protocols

How client/server communication happens

LucidDB, H2 -> PostgreSQL protocol

Sphinx -> MySQL protocol

Tokyo Cabinet / Tyrant -> memcached protocol

END*point*

# Verification

- 'memcapable' certifies memcached implementations

- Need automated, repeatable tests for complex systems (Cucumber?)

- More people connections between projects

END*point*

# Databases.
# Talking to each other.

# Thrift -> ThruDB
http://code.google.com/p/thrudb/

END*point*

# Thanks go to:

- Sheeri Cabral
- Josh Berkus
- Brian Aker
- Monty Taylor
- Stewart Smith
- Mark Atwood
- J Chris Anderson
- Jan Lehnardt
- Rick Hillegas
- Salvatore Sanfilippo

- Martin Kersten
- Robin Schumacher
- Vadim Tkachenko
- Justin Sheehy
- Nicholas Goodman, John Sichi, Joseph A. di Paolantonio
- Jay Pipes
- Tobias Downer
- Thomas Mueller
- Scott Deckelmann

LCA 2010

END *point*

# Questions?

END *point*

This work by Selena Deckelmann is licensed under a <u>Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States License</u>.

END*point*