# Sphinx

## Overview and quickstart

Vladimir Fedorkov // Sphinx Technologies Inc.
4Developers, 2011

# The Talk

- What Sphinx can do for you
- From basics to search cluster
- Upcoming release 2.0

# The Sphinx

- Free open source search server
- Begins 10 years ago as a full text daemon
- Now powerful, fast, relevant, scalable search engine.
- Not just full text search!

# Sphinx installations serves

- Over 16,000,000,000 (yes billions) documents
  - boardreader.com, over 5Tb data on about 40 boxes
- Over 200,000,000 queries/day (craigslist.org)
  - 2,000 QPS against 15 Sphinx boxes
- Also powers NetLog, Meetup, Slashdot, WikiMapia, and a few thousand other sites.

# Why Sphinx?

- ## 10-1000x vs MySQL on full-text searches
  - MySQL only behaves when indexes are in RAM

- ## 2-3x vs MySQL on non-full-text scans
  - On a single core
  - Because of less overheads

- ## 2-4x faster than Lucene on indexing and 2x faster on searches
  - Our (aged) internal benchmarks

# Sphinx
# Extended Full-text query syntax

- And, Or
  - hello | world, hello & world
- Not
  - hello -world
- Per-field search
  - @title hello @body world
- Search within first N symbols
  - @body[50] hello
- Phrase search
  - "hello world"
- Proximity search
  - "hello world"~10
- Distance support
  - hello NEAR/10 world
- Quorum matching and strict order support
- Custom weighting
- …and more

# Sphinx
# Not only Full-text

- Geo distance search
- MVA (i.e. page tags or multiple categories)
- UNIX timestamps
- Floating point values
- Strings
- Integers
- All the above combined altogether with FT search in a single query.

# Typical usages

- Item search
- Forum/blog posts search
- "Similar items/pages" service
- Misspelling correction service
  - Included in distribution
- Dating websites
  - Because of fast in-memory lookups

# Easy!

- Working out of the box
- You can run it on various platforms
  - Even AIX and iPhone!
  - PHP, Python, Java, Ruby, C binary protocol APIs are officially available.
  - .NET, Thinking Sphinx (for Rails) and few more available as third party plugins
- Can pull data from MySQL, PostgreSQL, MSSQL, any ODBC source and via XML pipe
- And more…

4Developers, Warsaw, 2011

# Uses MySQL protocol for SphinxQL!

```
$ mysql -h 0 -P 9306
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 1
Server version: 1.11-dev (r2569)

Type 'help;' or '\h' for help. Type '\c' to clear the current
   input statement.

mysql> SELECT * FROM lj1m WHERE MATCH('Sphinx') ORDER BY ts
   DESC LIMIT 3;
+---------+--------+------------+------------+
| id      | weight | channel_id | ts         |
+---------+--------+------------+------------+
| 7333394 |   1649 |     384139 | 1113235736 |
| 7138085 |   1649 |     402659 | 1113190323 |
| 7051055 |   1649 |     412502 | 1113163490 |
+---------+--------+------------+------------+
3 rows in set (0.00 sec)
```

4Developers, Warsaw, 2011

# How can I do that?

- Download
  - from http://www.sphinxsearch.com/download
- Install
  - from package or sources
- Configure
  - Define where and how to get data (configure data sources)
  - Tell Sphinx how to process data and how to search (configure indexes)
- Perform Indexing by running indexer
- Run searchd
- Bingo!

4Developers, Warsaw, 2011

# Configuration

```
source lj_source
{
    …
    sql_query = SELECT id, channel_id, ts, title, content FROM ljposts
    sql_attr_uint = channel_id
    sql_attr_timestamp = ts
    …
}

index lj
{
    source = lj_source
    path = /my/index/store/lj_index
}
```

# Indexing

```
$ ./indexer lj1m
Sphinx 1.11-dev (r2569)
Copyright (c) 2001-2010, Andrew Aksyonoff
Copyright (c) 2008-2010, Sphinx Technologies Inc (http://sph...

using config file './sphinx.conf'...
indexing index 'lj1m'...
collected 999944 docs, 1318.1 MB
sorted 224.2 Mhits, 100.0% done
total 999944 docs, 1318101119 bytes
total 158.080 sec, 8338160 bytes/sec, 6325.53 docs/sec
total 33 reads, 4.671 sec, 17032.9 kb/call avg, 141.5 msec/call
total 361 writes, 20.889 sec, 3566.1 kb/call avg, 57.8 msec/call
```

# Running searchd

```
$ ../bin/searchd -c sphinx.conf
Sphinx 1.11-dev (r2569)
Copyright (c) 2001-2010, Andrew Aksyonoff
Copyright (c) 2008-2010, Sphinx Technologies
   Inc (http://sphinxsearch.com)

using config file 'sphinx.conf'...
listening on 127.0.0.1:9312
listening on 127.0.0.1:9306
precaching index 'lj1m'
precached 1 indexes in 0.028 sec
```

# Sphinx

# How do I search from there?

- Sphinx API

```php
<?php
require ( "sphinxapi.php" );
$cl = new SphinxClient ();
$res = $cl->Query ( "my first query",
"my_index" );
var_dump ( $res );
// wham, bam, searching kinda done
```

# And then how do I search?

- SphinxSE

```
SELECT *
FROM sphinxsetable s
JOIN products p ON p.id=s.id
WHERE s.query='@title ipod'
ORDER BY p.price ASC

// or better!
... WHERE s.query='@title ipod;sort=attr_asc:price';
```

# And finally!

```
$ mysql -h 0 -P 9306
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 1
Server version: 1.11-dev (r2569)

Type 'help;' or '\h' for help. Type '\c' to clear the current
    input statement.

mysql> SELECT * FROM lj1m WHERE MATCH('Sphinx') ORDER BY ts DESC
    LIMIT 3;
+---------+--------+------------+------------+
| id      | weight | channel_id | ts         |
+---------+--------+------------+------------+
| 7333394 |   1649 |     384139 | 1113235736 |
| 7138085 |   1649 |     402659 | 1113190323 |
| 7051055 |   1649 |     412502 | 1113163490 |
+---------+--------+------------+------------+
3 rows in set (0.00 sec)
```

4Developers, Warsaw, 2011

# SphinxQL

Our own implementation of MySQL protocol

- Our own SQL parser
- **MySQL not required**!
- Any **client** library (eg. PHP's or .NET) should suffice
- All new features will initially appear in SphinxQL

# RT indexes

- Push model instead of Pull for on-disk indexes
  - via INSERT/UPDATE/DELETE
- Update data on the fly
- Formally "soft-realtime"
  - As in, most of the writes are very quick
  - But, not guaranteed to complete in fixed time
- Transparent for application

# RT indexes, the differences

- Indexing is SphinxQL only
  - mysql_connect() to Sphinx instead of MySQL
  - mysql_query() and do INSERT/REPLACE/DELETE as usual
- Searching is transparent
  - SphinxAPI / SphinxSE / SphinxQL all work
  - We now prefer SELECT that we have SphinxQL :)
- Some features are not yet (!) supported
  - MVA, geosearch, prefix and infix indexing support to be implemented

# Scale? Scale!

- Utilize multicore servers
- Spread load across several boxes
- Shard the data

# Scale? Scale!

- Create several local indexes
- Create distributed index and query
  - local indexes from the same box
  - remote Sphinx instances

# Scaling part one: data sources

```
source lj_source
{
    …
    sql_query          = SELECT id, channel_id, ts, title,
    content FROM ljposts WHERE id>=$start and id<=$end
    sql_query_range    = SELECT 1, 7765020
    sql_attr_uint      = channel_id
    sql_attr_timestamp = ts
    …
}


source lj_source2 : lj_source
{
     sql_query_range  = SELECT 7765020, 10425075
}
```

# Scaling part two: local indexes

```
index ondisk_index1
{
  source            = source1
  path              = /path/to/ondisk_index1
  stopwords         = stopwords.txt
  charset_type      = utf-8
}

index ondisk_index2 : ondisk_index1
{
  source            = source2
  path              = /path/to/ondisk_index2
}
```

# Scaling part three: distributed indexes

```
index my_distribited_index1
{
  type       = distributed
  local      = ondisk_index1
  local      = ondisk_index2
  local      = ondisk_index3
  local      = ondisk_index4
}
…
  dist_threads = 4
…
```

# Scaling part three: distributed indexes

```
index my_distribited_index2
{
   type     = distributed
   agent    = 192.168.100.51:9312:ondisk_index1
   agent    = 192.168.100.52:9312:ondisk_index2
   agent    = 192.168.100.53:9312:rt_index
}
```

# Sphinx

# Beyond the basics

- Tokenizing settings
- Wordforms support
- 1-grams
- HTML processing
- SQL and IO throttling
- Arbitrary expressions
- Prefix/infix indexing
- Blended characters
- Hitless indexing
- Different rankers
- Some more features…

4Developers, Warsaw, 2011

# Upcoming 2.0.1 release

- Improved SphinxQL
- dict=keywords
- Zones, sentences, paragraphs support
- Multi-threaded snippet batches support
- UDF support (CREATE/DROP FUNCTION)
- Support for ORDER BY, GROUP BY, WITHING GROUP ORDER BY for strings
- New query log file format
- 35 more new features

# Sphinx today

## We hiring!

Consultants, support engineers,
Q/A engineer and technical writer wanted!

http://sphinxsearch.com/about/careers/

Just let me know or
mail us at job2011@sphinxsearch.com

4Developers, Warsaw, 2011

# Questions?



http://sphinxsearch.com