# Solr for newbies

https://hectorcorrea.com/solr-for-newbies

Hector Correa
hector_correa@princeton.edu
Princeton University

code{4}lib

# Useful Links

**Workshop links**

http://hectorcorrea.com/solr-for-newbies

**Code4Lib code of conduct**

https://2023.code4lib.org/conduct/

# Workshop Outline

**1. Introduction**
(concepts, quick tour, installation)

**2. Schema**
(fields, field types, query/ index analyzers, tokenizers)

**3. Searching**
(query parsers, search params, facets, highlighting)

**4. Miscellaneous**
(directories, configuration, synonyms, spellcheck)

**1. Introduction**
(concepts, quick tour, installation)

**2. Schema**
(fields, field types, query/ index analyzers, tokenizers)

**3. Searching**
(query parsers, search params, facets, highlighting)

**4. Miscellaneous**
(directories, configuration, synonyms, spellcheck)
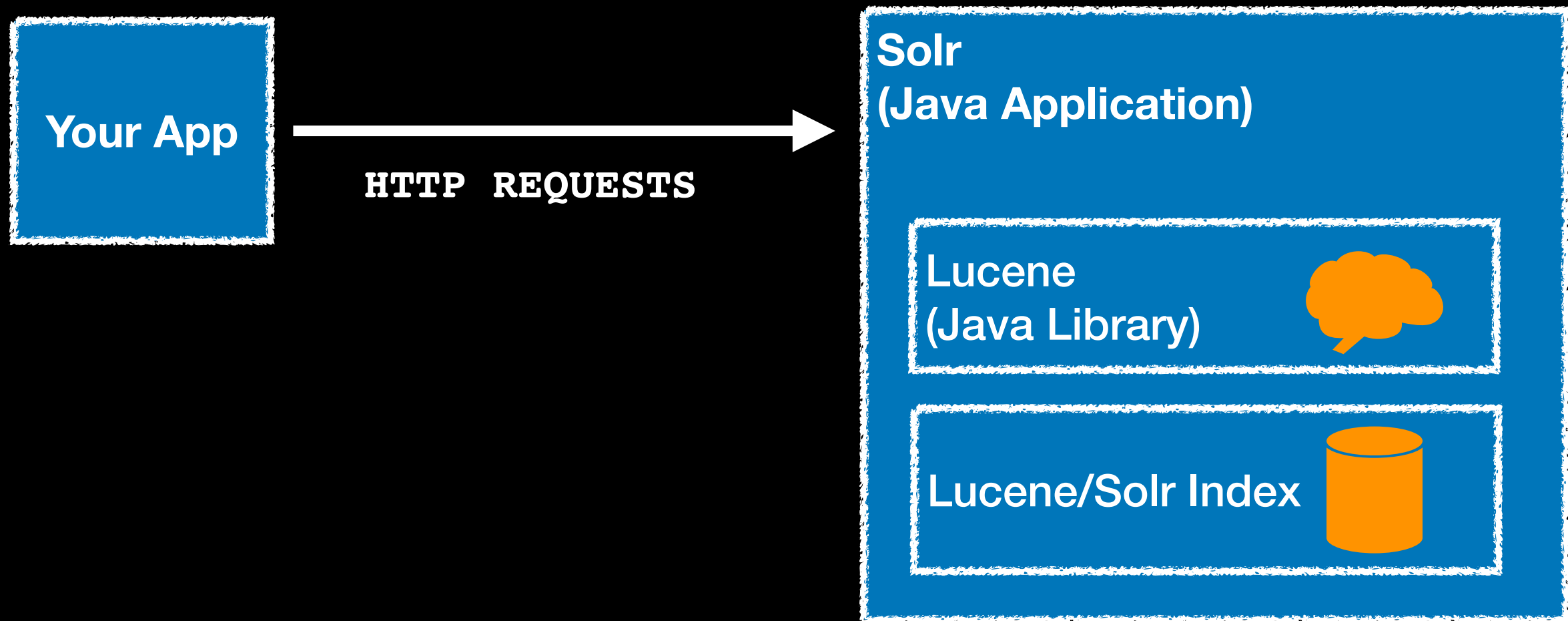
# What is Solr

"Solr is the popular, blazing-fast,
open source enterprise **search platform**
built on Apache Lucene."
- Solr's Home Page

"Solr is a scalable, ready-to-deploy
enterprise **search engine** that's optimized
to search large volumes of text-centric data
and return results sorted by relevance."
- Solr in Action [p. 4]

# What is Solr

"Solr is the popular, blazing-fast,
open source enterprise **search platform**
built on Apache Lucene."
- Solr's Home Page

"Solr is a scalable, ready-to-deploy
enterprise **search engine** that's optimized
to search large volumes of **text-centric data**
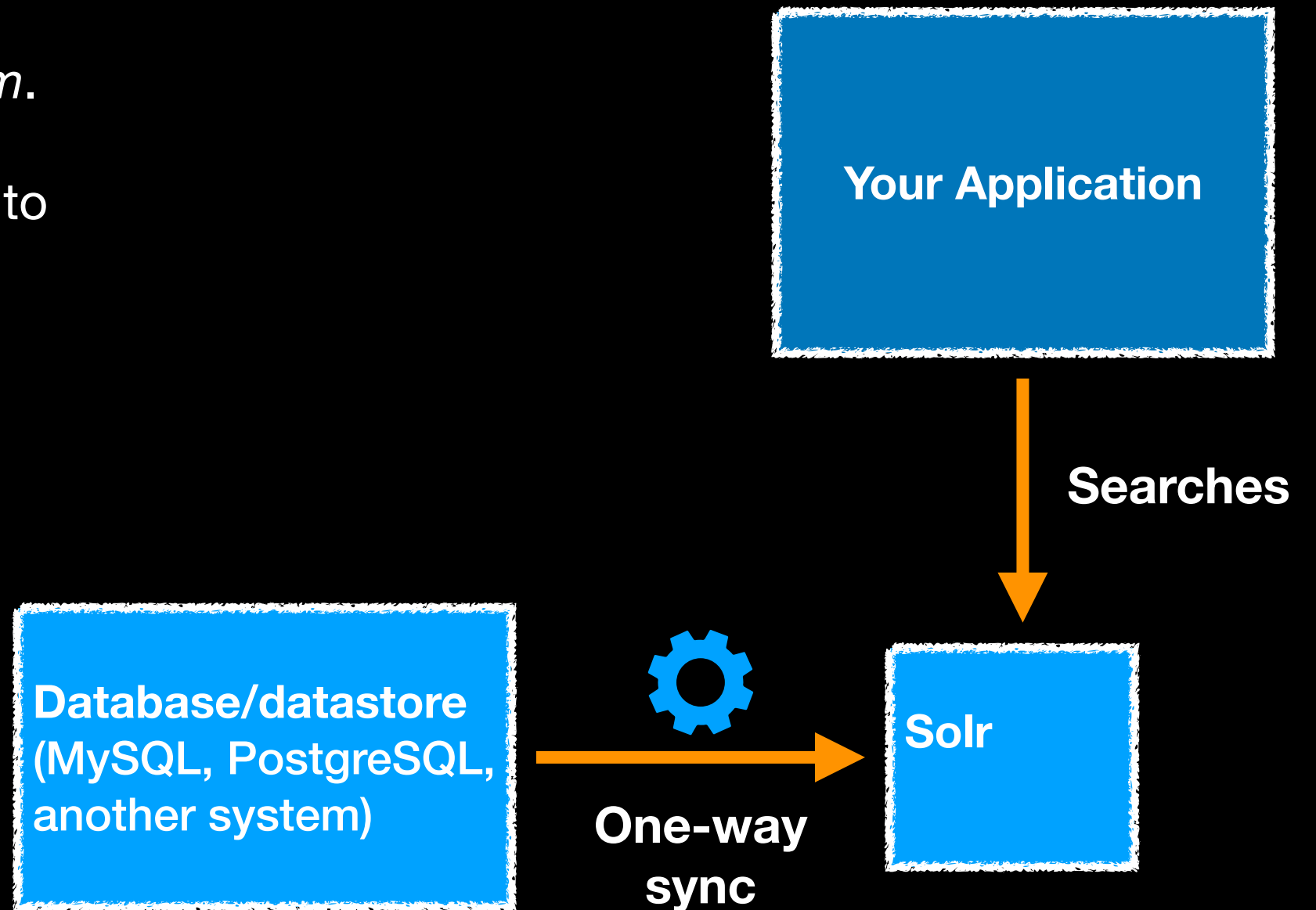and return **results sorted by relevance**."
- Solr in Action [p. 4]

# Your App, Solr, and Lucene

**Your App** → **HTTP REQUESTS** → **Solr (Java Application)**

**Lucene (Java Library)**

**Lucene/Solr Index**

# Typical Architectures I

Your application searches via Solr, *but the data is maintained in another system*.
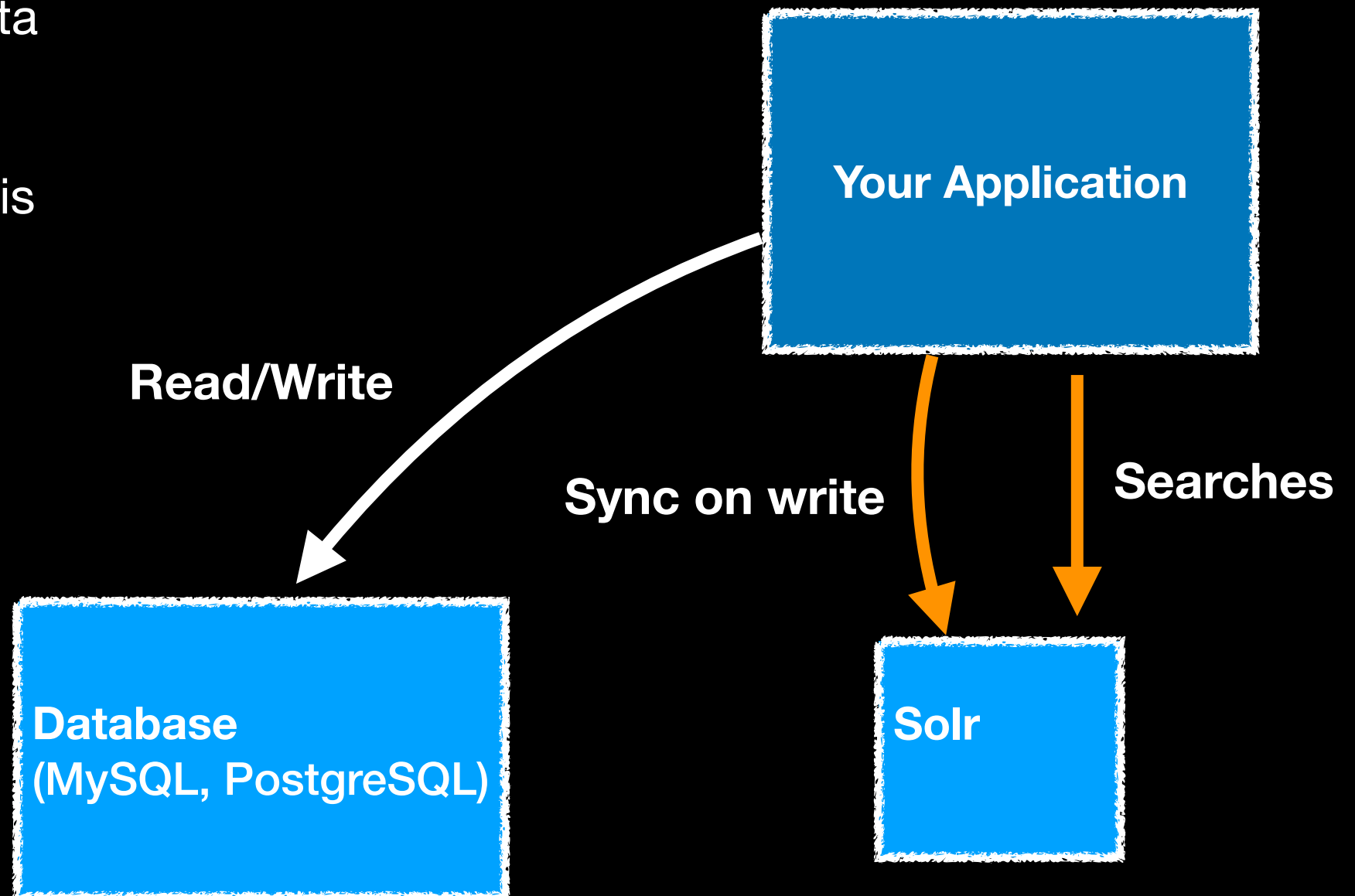
Blacklight applications tend to follow this pattern.

**Your Application**

**Searches**

**Database/datastore (MySQL, PostgreSQL, another system)**

**One-way sync**

**Solr**

# Typical Architectures II

Your application uses a database to maintain the data and Solr for searches.
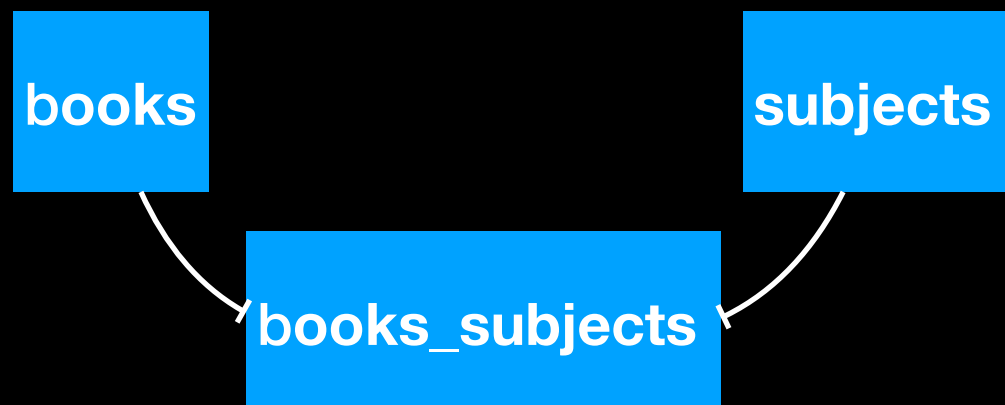
VIVO and SamVera follow this pattern.

**Your Application**

**Read/Write**

**Sync on write**

**Searches**

**Database**
**(MySQL, PostgreSQL)**

**Solr**

# Document Model
## (how Solr *stores* your data)

| id | book_title | subjects |
|----|-----------|----------|
| 1 | Princeton guide for dog owners | animals, guides |
| 2 | Princeton tour guide | guides |
| 3 | Cats and dogs | animals |

# Relational Model

books

subjects

books_subjects

# Document Model

```
solr_doc: {
  id:"1",
  title:"Princeton guide for dog owners",
  subjects: ["animals", "guides"]
}
```

# Solr Documents are flat*

```
your_data:
{
  id:"9041",
  title:"Using Qualitative Inquiry to Promote…",
  authors: [
    {uri:"http://somebody/51", name: "Loya, Karla"},
    {uri:"http://somebody/82", name: "Kimball, Ezekiel"}
  ],
  subjects: ["higher education", "org theory"]
}
```

*data in Solr is flatten*

```
    solr_doc: {
      id:"9041",
      title:"Using Qualitative Inquiry to Promote…",
      authors_uri: ["http://somebody/51", "http://somebody/82"],
      authors_name: ["Kimball, Ezekiel", "Loya, Karla"],
      subjects: ["higher education", "org theory"]
    }
```

# Inverted Index
## (how Solr *indexes* your data)

| id | title | subjects |
|---|---|---|
| 1 | Princeton guide for dog owners | animals, guides |
| 2 | Princeton tour guide | guides |
| 3 | Cats and dogs | animals |

## Traditional Index

| id | title |
|---|---|
| 3 | Cats and dogs |
| 1 | Princeton guide for dog owners |
| 2 | Princeton tour guide |

## Inverted Index

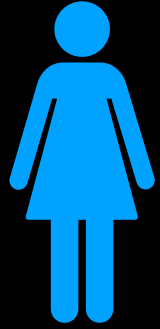| key | ids |
|---|---|
| princeton | 1, 2 |
| owners | 1 |
| dogs | 1, 3 |
| guide | 1, 2 |
| tour | 2 |
| cats | 3 |

**1. Introduction**
(concepts, quick tour, installation)

**2. Schema**
(fields, field types, query/ index analyzers, tokenizers)

**3. Searching**
(query parsers, search params, facets, highlighting)

**4. Miscellaneous**
(directories, configuration, synonyms, spellcheck)

# Adding a document to Solr

**HTTP POST**
http://localhost/solr/**bibdata**/**update**

```
{
 id:"1",
 title_txt_en:"history of medicine",
 subject_s: "medicine",
 abstract_txt: "this book is about..."
}
```

**Solr**

**bibdata** core

**/update Handler**
(solrconfig.xml)

**Index** Analyzers
tokenizer + filters for each field
(schema.xml)

Lucene Index

# Adding a document to Solr

## Your data

```
{
  id:"1",
  title_txt_en:"history of medicine",
  subject_s: "medicine",
  abstract_txt_en: "this book is about..."
}
```

## + Solr's Schema

```
<field name="id" type="string" multiValued="false" />

<dynamicField name="*_s" type="string" />
<dynamicField name="*_txt_en" type="text_en" />
<dynamicField name="*_txt" type="text_general" />
```

## Gives

id and subject will be handled as a string
title and abstract will be handled as text_en

# Adding a document to Solr

id and subject will be handled as a string

```
$ curl localhost:8983/solr/bibdata/schema/fieldtypes/string

"fieldType":{
    "name":"string",
    "class":"solr.StrField",
    "sortMissingLast":true,
    "docValues":true
}
```

# Adding a document to Solr

Title and abstract will be handled as a text_en

```
$ curl localhost:8983/solr/bibdata/schema/fieldtypes/text_en

"fieldType":{
    "name":"text_en",
    "class":"solr.TextField",
    "positionIncrementGap":"100",
    "multiValued":true,
    "indexAnalyzer":{
      "tokenizer":{StandardTokenizerFactory},
      "filters":[StopFilterFactory, LowerCaseFilterFactory,
EnglishPossessiveFilterFactory, PorterStemFilterFactory}]},
    "queryAnalyzer":{
      "tokenizer":…,
      "filters":[…]}}
}
```
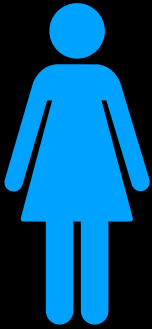
# Workshop Outline

**1. Introduction**
(concepts, quick tour, installation)

**2. Schema**
(fields, field types, query/index analyzers)

**3. Searching**
(query parsers, search params, facets, highlighting)

**4. Miscellaneous**
(directories, configuration, synonyms, spellcheck)

# Searching for documents in Solr

**HTTP GET**
http://localhost/solr/**bibdata**/**select**
?q=subject:medicine

Documents
Facets
Highlighting
…

**Solr**

**bibdata** core

/select Handler
(solrconfig.xml)

Query Parser (e.g. eDisMax)

**Query Analyzers**
tokenizer + filters for each field
(schema.xml)

Lucene Index

# Thanks and good luck

Stay in touch

hector_correa@princeton.edu
`https://hectorcorrea.com/solr-for-newbies`