# Apache Solr

DAPI . Information Description, Storage and Retrieval Course
MIEIC, 2020/21 Edition

Sérgio Nunes
DEI, FEUP, U.Porto

*Work in progress*

# Plan for Today

➔ Questions?

➔ Groups Presentations (~90 min)

➔ Break

➔ Milestone #2 Overview

➔ Solr overview
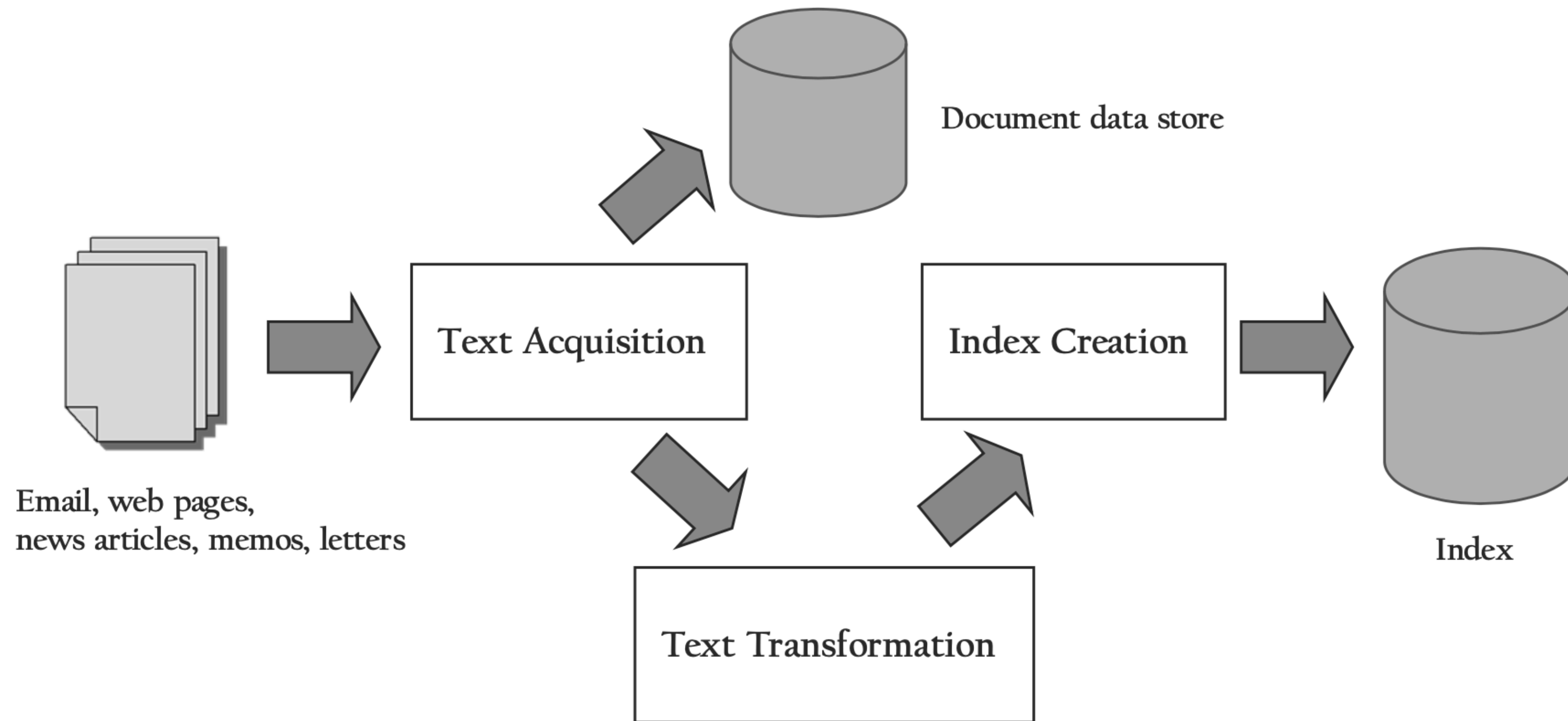
# Milestone #2 — Information Retrieval

# Milestone #2

➔ Goal: index dataset to support querying using free-text

➔ Use open-source tools (i.e. Solr); decide on the document granularity; decide on the search filters.

➔ Expected actions:

  ➔ Choose the information retrieval tool (Solr, Lucene, Terrier, Elasticsearch, …);

  ➔ Analyze the documents and identify their indexable components;

  ➔ Identify search parameters that will be offered to the users;

  ➔ Use the tool API to generate indexes;

  ➔ Use the tool API to configure the answer to queries;

  ➔ Demonstrate the indexing and retrieval processes;

  ➔ Evaluate the results, (ideally) comparing different ranking formulas.

➔ More information at https://web.fe.up.pt/~ssn/dokuwiki/teach/dapi/202021/delivery2/index

# Search Engine Overview

# Architecture of a Search Engine

➔ Two primary goals of a search engine:

    ➔ *effectiveness* (quality) — retrieve the most relevant set of documents;

    ➔ *efficiency* (speed) — present the results as quickly as possible;

➔ Search engines are architected to two support two major functions:

    ➔ *indexing process* — build the structures to enable search;

    ➔ *querying process* —use the structures to produce a ranking;

# The Indexing Process



Document data store

Email, web pages,
news articles, memos, letters

Text Acquisition

Text Transformation

Index Creation

Index

# Blocks of the Indexing Process

➜ Text Acquisition

➜ crawler; conversion; document data store.

➜ Text Transformation

➜ parser; stopping; stemming; link extraction; information extraction; classifier.

➜ Index Creation

➜ document statistics; weighting; inversion; index distribution;

# The Querying Process

# Blocks of the Querying Process

➜ User Interaction

   ➜ query input; query transformation; results output.

➜ Ranking

   ➜ scoring; performance optimization; distribution.

➜ Evaluation

   ➜ logging; ranking analysis; performance analysis.

# Apache Solr

# Apache Solr

➔ Solr is a search server built on top of Apache Lucene, an open source, Java-based, information retrieval library. Standard steps:

  ➔ Define the schema, to tell Solr about the contents of documents it will be indexing;

  ➔ Feed Solr documents for which your users will search;

  ➔ Expose search functionality in your application.

➔ Solr offers support for the simplest keyword searching through to complex queries on multiple fields and faceted search results.

➔ Because Solr is based on open standards, it is highly extensible. Solr queries are simple HTTP request URLs and the response is a structured document: mainly JSON, but it could also be XML, CSV, or other formats.

# Example Solr Integration

# Solr Features



From: *Trey Grainger and Timothy Potter. Solr in Action, Manning Publications, 2014.*

# The Inverted Index (again)



From: *Trey Grainger and Timothy Potter. Solr in Action, Manning Publications, 2014.*

# Solr Command Line Tool

# Solr Admin Console

# Tasks

➔ Finish and submit Milestone #1 report.

➔ Review goals and organize work for Milestone #2.

➔ Experiment with full-text indexing tools.

   ➔ Apache Solr Tutorial — https://lucene.apache.org/solr/guide/solr-tutorial.html

   ➔ Experiment with other collections (e.g. project, personal documents, etc).

➔ Anticipate indexing and search tasks on the working dataset.

➔ **Next week**: finish and submit Milestone #1 report.

# References

➜ Apache. Solr Tutorial. https://lucene.apache.org/solr/guide/solr-tutorial.html

➜ Apache Solr Reference Guide. https://lucene.apache.org/solr/guide/

➜ Trey Grainger and Timothy Potter. Solr in Action, Manning Publications, 2014.

➜ W. Bruce Croft, Donald Metzler, Trevor Strohman, Search Engines: Information Retrieval in Practice, Pearson, 2009. http://ciir.cs.umass.edu/downloads/SEIRiP.pdf