# SANSA
# in the Palm of your Browser

Carsten Felix Draschner - SDA Research

# Overview

- get Databricks
- setup Databricks
    - get SANSA Jar
    - get sample files
    - create cluster
    - set spark enviroment
    - use SANSA jar
    - create notebook
    - use sansa modules

# Databricks

- Similar concepts
  - jupyter notebooks
  - google colab
  - Scala
  - Apache Spark
  - AWS
- Advantages
  - No istallation required
  - available over browser
  - native scala spark notebooks
  - no local performance needed

# Databricks Registration

- Databricks FAQ
  - https://databricks.com/de/product/faq/community-edition
- Login
  - https://community.cloud.databricks.com/login.html
- Or create and use for free
  - https://databricks.com/try-databricks
  - 15GB Ram, 2 Core Cluster SAmple Cluster

# Setup Databricks - Upload needed Data - jar

- Most recent SANSA Release Jar Available on Github Page:
  - https://github.com/SANSA-Stack/SANSA-Stack/releases
- Or through this link directly: 234mb fat jar
  - sansa-stack-spark_2.12-0.8.0-RC1-jar-with-dependencies.jar

# Setup Databricks - Upload needed Data - jar

# Setup Databricks - Upload needed Data - jar

# Setup Databricks - Upload needed Data - jar

# Setup Databricks - Upload needed Data - jar



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Databricks - Upload needed Data - jar

# Setup Databricks - Upload needed Data - data

- Sample data
  - LMDB
    - http://www.cs.toronto.edu/~oktie/linkedmdb/linkedmdb-18-05-2009-dump.nt
  - Sample Data
    - E4toE5.zip

# Setup Databricks - Upload needed Data - data



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Databricks - Upload needed Data - data



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Databricks - Upload needed Data - data



SANSA in the Palm of your Browser - Carsten Felix Draschner
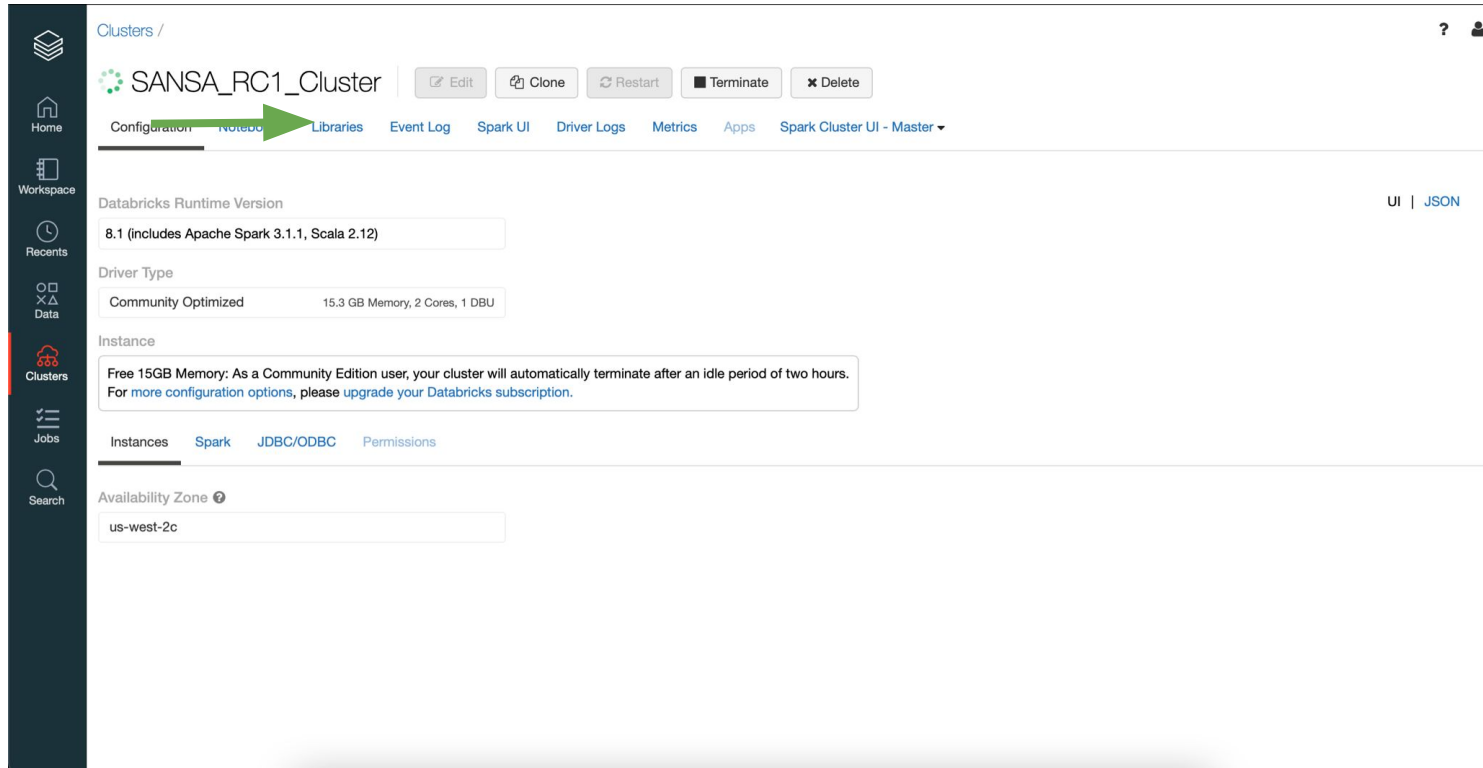
# Setup Cluster

- set spark (each one line)

  - `spark.databricks.delta.preview.enabled true`

  - `spark.serializer org.apache.spark.serializer.KryoSerializer`

  - `spark.kryo.registrator net.sansa_stack.rdf.spark.io.JenaKryoRegistrator,`
    `net.sansa_stack.query.spark.sparqlify.KryoRegistratorSparqlify`

- set jar

# Setup Cluster



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - Give it a Name

# Setup Cluster - Set Spark

# Setup Cluster - Set Spark

# Setup Cluster - Set Spark

# Setup Cluster - Set Spark

# Setup Cluster - Set Spark



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - Create Cluster



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Setup Cluster - set SANSA JAR

# Setup Cluster - set SANSA JAR

# Setup Cluster - set SANSA JAR

# Import Notebook

- Notebook
  - [https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6924783690087984/3044317244801485/8524188481975304/latest.html](https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6924783690087984/3044317244801485/8524188481975304/latest.html)

# Import Notebook

SANSA in the Palm of your Browser - Carsten Felix Draschner

# Import Notebook

# Import Notebook

# Import Notebook

# Run Notebook - Select Cluster

SANSA in the Palm of your Browser - Carsten Felix Draschner

# Run Notebook - Select Cluster

# Run Notebook - Run all cells

# You Made It!!!

# Create Notebook

- Create Notebook
- Attach Cluster
- Read In Data
- Perform DistSim Modules

# Create Notebook - Create Notebook



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Create Notebook - Create Notebook

# Create Notebook - Create Notebook



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Create Notebook - Specify Cluster

# Create Notebook - Specify Cluster



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Code in Notebook - Imports

```
import net.sansa_stack.ml.spark.utils.{FeatureExtractorModel, SimilarityExperimentMetaGraphFactory}
import net.sansa_stack.rdf.spark.io._
import org.apache.jena.graph
import org.apache.jena.riot.Lang
import org.apache.jena.sys.JenaSystem
import org.apache.spark.ml.feature.{CountVectorizer, CountVectorizerModel}
import org.apache.spark.ml.linalg.Vector
import org.apache.spark.rdd.RDD
import org.apache.spark.sql.functions.{col, udf}
import org.apache.spark.sql.{DataFrame, Dataset}
```

# Code in Notebook - Imports



SANSA in the Palm of your Browser - Carsten Felix Draschner

# Code in Notebook - Imports

# Code in Notebook - Imports



SANSA in the Palm of your Browser - Carsten Felix Draschner