

Open platform of civil society organizations

Challenges, architecture, techniques and technologies in development of the most important system of IPEA

Bio

Technical coordinator (Researcher II) at IPEA

MBA Professor at INFNET

MSc student at PESC/UFRJ

Research interests: Data Mining & Machine Learning in Dynamic Environments,
Complex Networks, Big Data for Social Development

Summary

Introduction

Challenges

Architecture

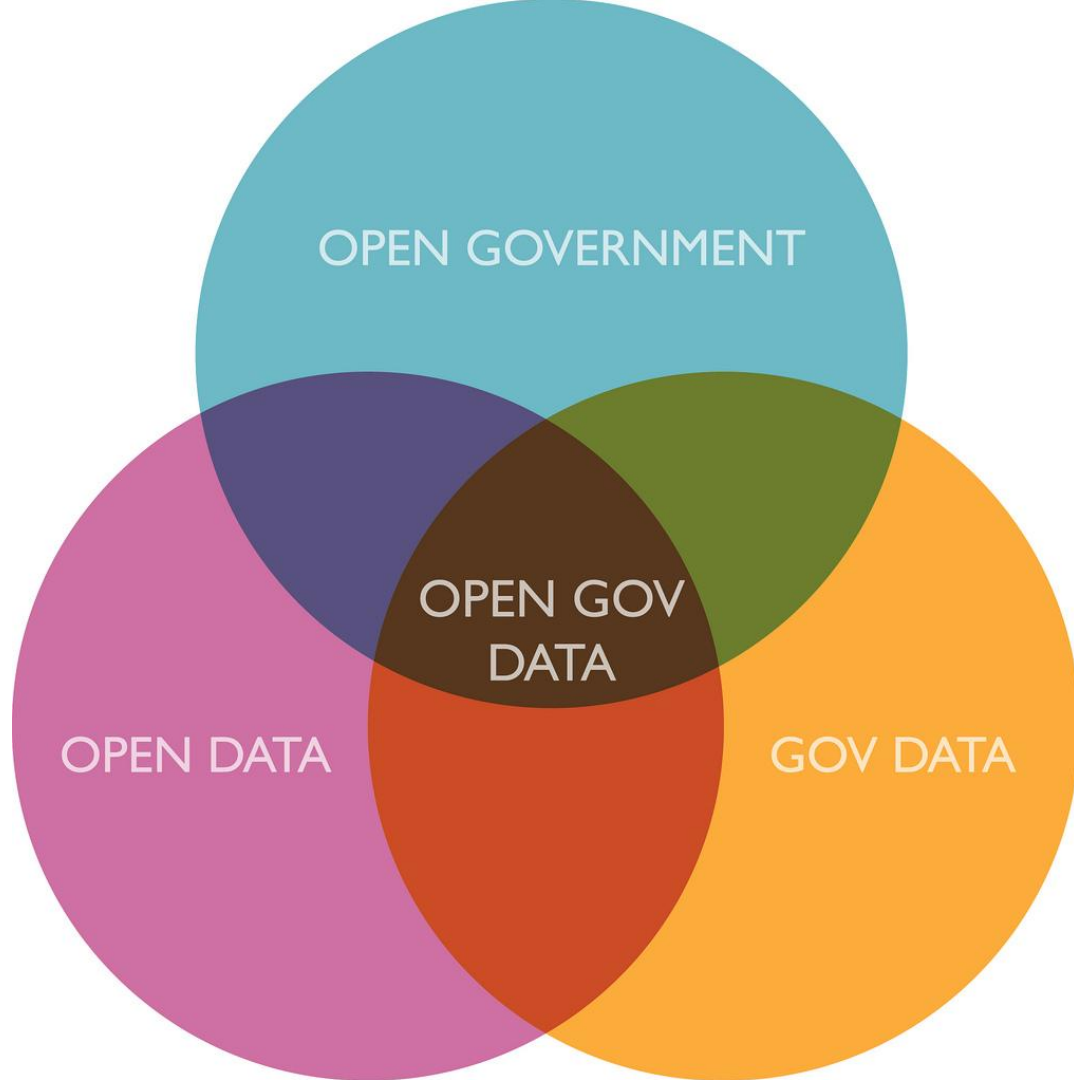
Techniques

Technologies

Current state

Future steps

Introduction



IPEA

Instituto de Pesquisa Econômica Aplicada (IPEA) is a traditional Brazilian federal public foundation linked to the Ministry of Planning, Development and Management

Focus on support formulation and reformulation of public policies and Brazilian development programs

Traditional areas: Economics, stats, sociology, public policy research

New areas: Information technology, big data management, computer science

Civil society organizations

Civil society organization is the "aggregate of non-governmental organizations and institutions that manifest interests and will of citizens."

Data about these organizations must to be open for society

Brazilian government is investing in open data and new ways to allow research in this data

One of the most expected systems at IPEA since this one will be the baseline to get info about all organizations of the country

Development team

Effort to build this system is divided in two teams located at Brasília and Rio

The project has around 15 members originated from many areas: Sociology, Stats, Computer Science and others

In the systems development team: DBA, designers, devops, developers and testers

Almost all members are graduated

Challenges

Big data

Process and visualize hundreds of thousands of points in the view:

- ~400 millions of registers

Integrate many different databases from other governmental domains

- The system is linked with 6 different big datasets

The ecosystem must be “Data Lake” ready

- The API is only a single independent application that will use a data from a big repository

BI, machine learning and data mining

Real time analytics

- Helping decision makers and researchers through the data

Automatic classification

- Organizations must be labeled and it means high costs to do it

Financial patterns

- Patterns helps to understand the growth of a NGO
- Fraud detection
- Financial linking between organizations

Open data concerns

Code and data must follow the open data plan

- Brazilian government demands that systems follows an open data plan

Needs strong security policies

- A governmental system must be ready to suffer cyber attacks

Data privacy

- Certain data attributes cannot be published or leaked

Continuous deployment

Get code from repository automatically

Automatic tests in front-end and back-end

Each system must be isolated

Build must be automatic

How achieve success using all things above using a simple solution but mature enough to use inside an big institution ?

Technologies

Databases and ETL process



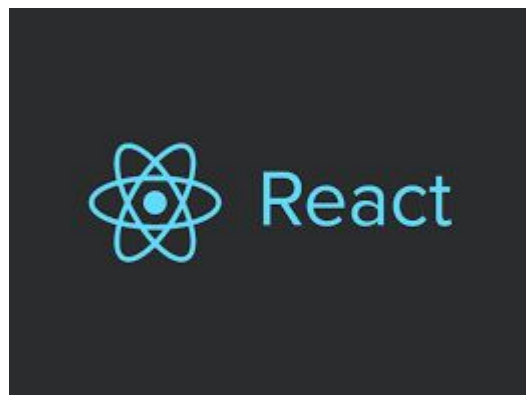
mongoDB®



Back-end



Front-end



Tests and automation



Jenkins



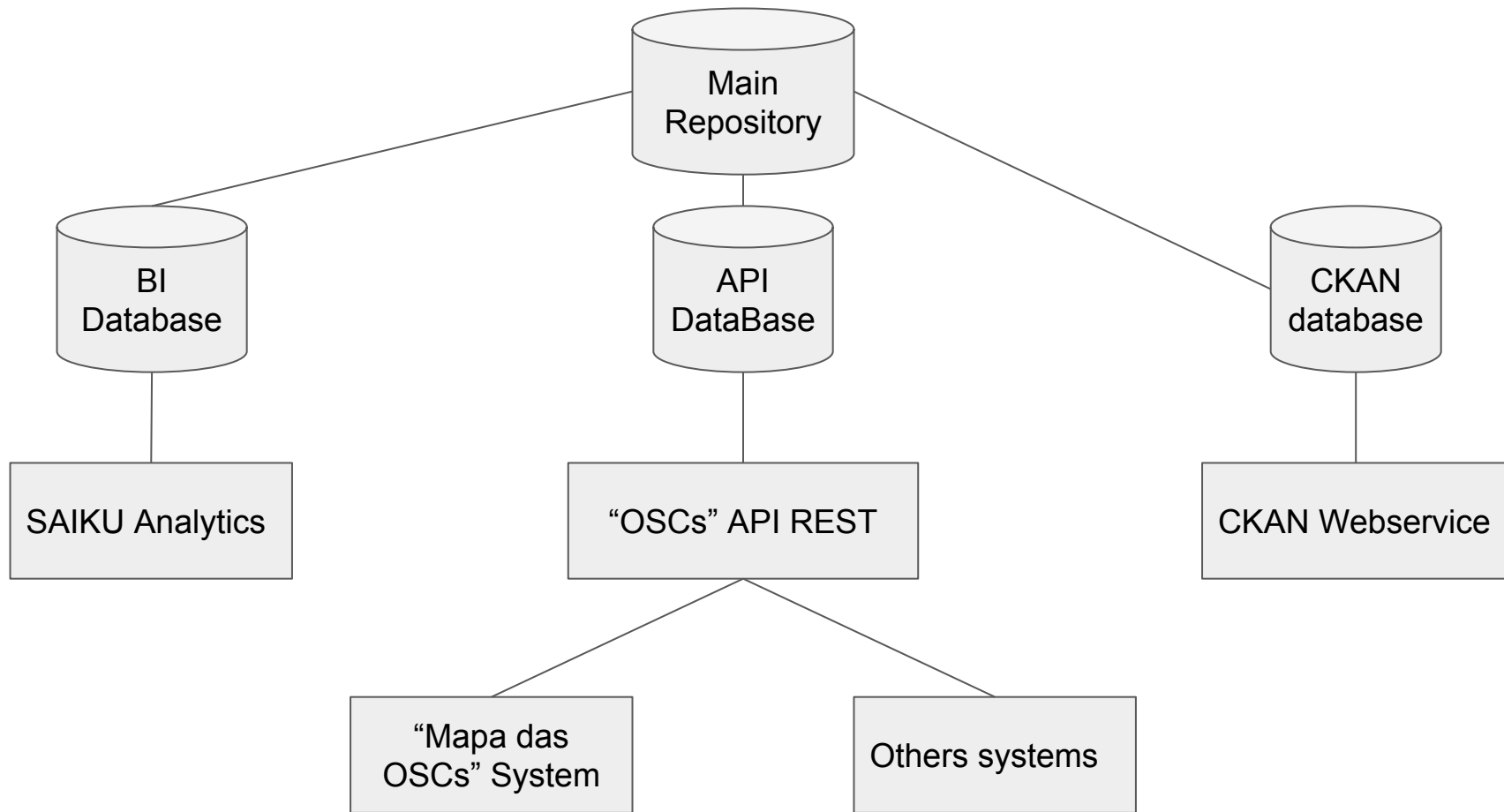
docker



Dependency Manager for PHP



Architecture



Database architecture

Main repository

- MongoDB
- Repository only for reading
- Raw data

BI database

- Postgres
- OLAP cube for analytics (SAIKU)

Database architecture

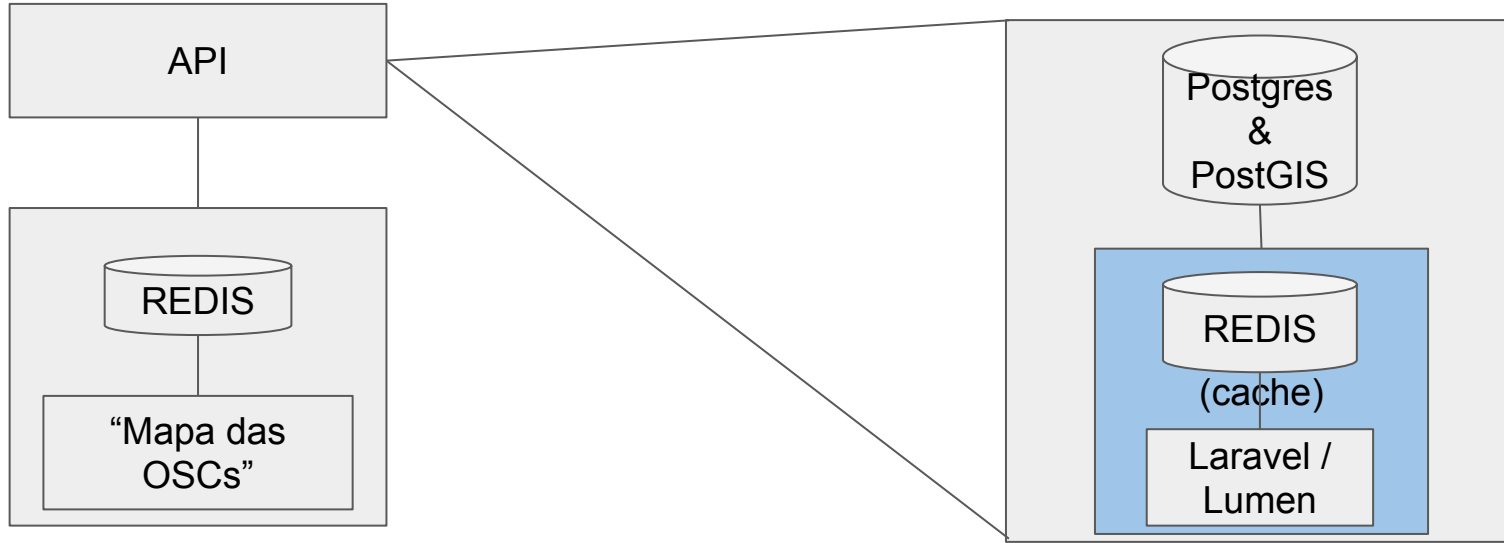
CKAN database

- Postgres
- Feeds the CKAN webservice

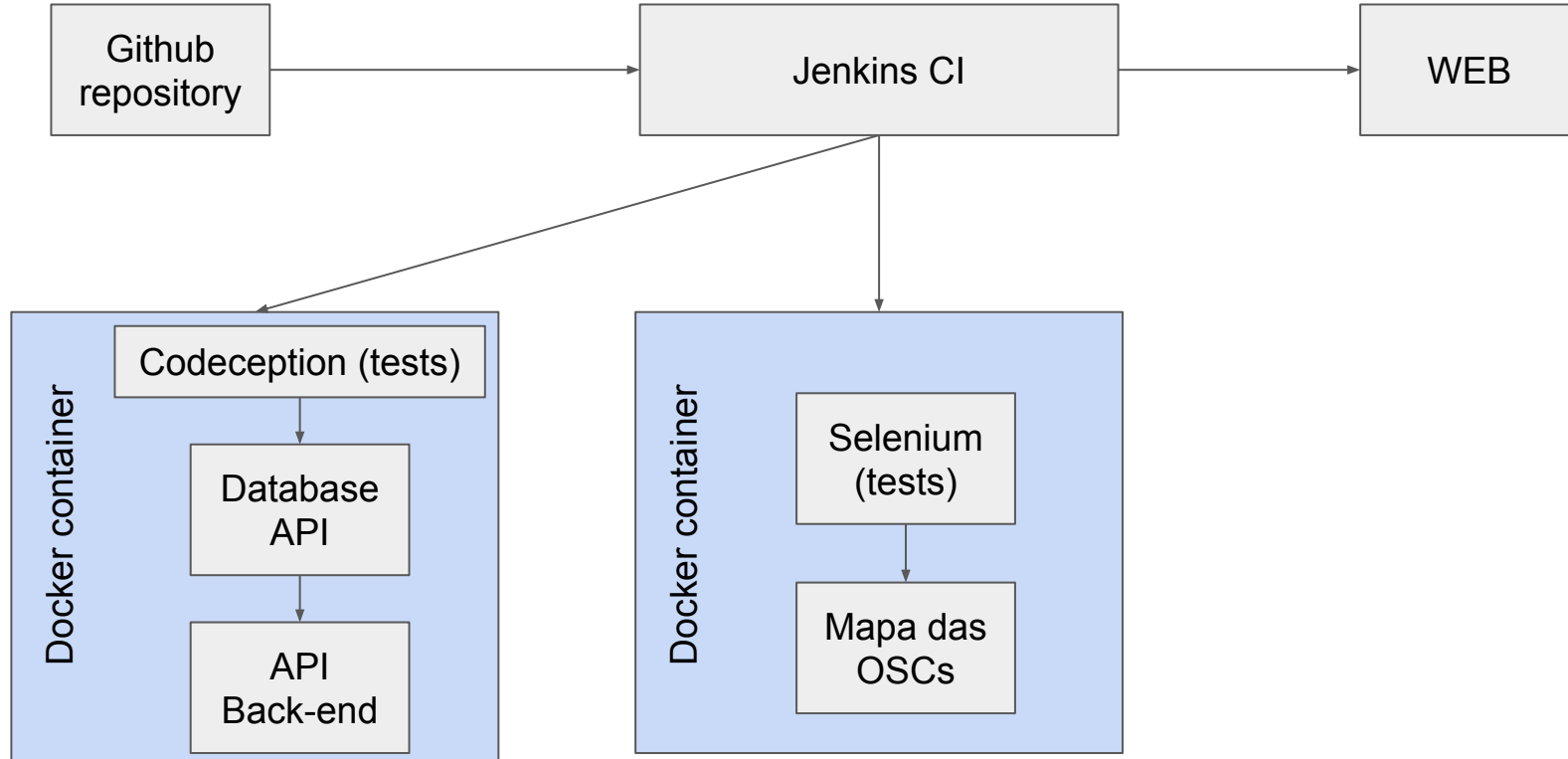
Database API

- Postgres & PostGIS
- Clean database for IPEA applications concern about CSOs

API architecture



Deployment



Techniques

Database

All SQL inside functions

R-Tree indexing at geometry bounds

Materialized views serving the API

Triggers updating views and cache

Databases are isolated from each other

Back-end

Data sent by chunks:

- Big chunks of data sent by PHP using functions `flush()`; and `sleep()`;

All queries pass through a middleware before to go the API:

- Search becomes a key and response becomes a value. These info are stored inside REDIS. Only new queries are sent to API.
- API also has his own cache (REDIS). Only new queries are sent to database.

Front-end

Transform almost everything in web components using React.js

```
var DropdownMenu = React.createClass({
  renderList: function(){
    var elems = [];
    for (var i=0; i< this.props.submenu.length; i++) {
      var l = this.props.submenu[i];
      //se titulo vir vazio entende-se que devemos adicionar um separador
      if (l.text=='') elems.push(<li className="divider"></li>);
      else elems.push(<li><a href={l.link}>{l.text}</a></li>);
    }
    return elems;
  },
  render: function () {
    return (<div>
      <a href="#" className="dropdown-toggle" data-toggle="dropdown" role="button" aria-expanded="false">{this.props.titulo}</a>
      <ul className="dropdown-menu" role="menu">{this.renderList()}</ul>
    </div>);
  }
});
```

Front-end

Require.js for transform
components in modules

```
require(['react', 'jsx!components/Util'], function (React) {

    require(['componenteBlocoDeTexto'], function(BlocoTexto){
        function BlocoDeTexto(titulo, formato){
            this.titulo = titulo;
            this.formato = formato;
        }

        var csv = "O formato CSV(Comma Separated Values) é um dos formatos mais utilizados para a troca de dados entre";
        var xls = "Formato padrão do Microsoft Excel, o XLS(XML Spreadsheet) tem uma qualidade razoável e é simples de";
        var xml = "O formato XML (eXtensible Markup Language) é popularmente utilizado para a transferência de dados binários";
        var json = "O formato JSON (JavaScript Object Notation) é um formato de transferência de dados que apresenta características próprias";
        var titulos = ["CSV", "XLS", "XML", "JSON"];
        var formatos = [csv, xls, xml, json];

        var blocosDeTexto = [];
        for (var i=0; i<titulos.length; i++){
            blocosDeTexto.push(new BlocoDeTexto(titulos[i], formatos[i]));
        }

        BlocoTexto = React.createFactory(BlocoTexto);
        ReactDOM.render(BlocoTexto({dados:blocosDeTexto}), document.getElementById("bloco_texto_formato_dados"));
    });

    require(['componenteDropdown'], function(Dropdown){
        var arquivosRetornados, arquivosEnviados;
        arquivosRetornados = arquivosEnviados = ["XML", "JSON", "CSV"];
        var periodicidade = ["Dia(s)", "Semana(s)", "Mês(es)"];

        Dropdown = React.createFactory(Dropdown);

        ReactDOM.render(Dropdown({list: arquivosRetornados}), document.getElementById("arquivo_retornado_dropdown"));
        ReactDOM.render(Dropdown({list:periodicidade}), document.getElementById("periodicidade_dropdown"));
        ReactDOM.render(Dropdown({list:arquivosEnviados}), document.getElementById("tipo_arquivo_dropdown"));
    });
});
```

Front-end

Dependencies are managed
with require.js config file

```
require.config({
  baseUrl: "js/",
  paths: {
    "react": "libs/react-15.3.1/react-with-addons.min",
    "jsx": "libs/jsxcompiler/jsx",
    "text": "libs/require-2.3.2/text",
    "JSXTransformer": "libs/jsxcompiler/JSXTransformer",
    "babel": "libs/babel-core/5.8.24/browser.min",
    "jquery": "libs/jquery-3.1.0/jquery-3.1.0.min",
    "jquery-ui": "libs/jquery-ui-1.12.0/jquery-ui",
    "bootstrap": "libs/bootstrap-3.3.7/bootstrap.min",
    "d3": "libs/nv-d3/d3.v3",
    "nv.d3": "libs/nv-d3/nv.d3",
    "nv.d3.lib": "libs/nv-d3/nv.d3.lib",
    "stream": "libs/nv-d3/stream-layers",
    "tablesaw": "libs/tablesaw-3.0/tablesaw",
    "tablesaw-init": "libs/tablesaw-3.0/tablesaw-init",
    "datatables.net": "libs/DataTables/DataTables-1.10.12/js/jquery.dataTables.min",
    "datatables-responsive": "libs/DataTables/Responsive-2.1.0/js/dataTables.responsive.min",
    "leaflet": "libs/leaflet-0.7.7/leaflet",
    "leafletCluster": "libs/leaflet-0.7.7/cluster.min",
    "google": "libs/google",
    "rotas": "rotas"
  },
  shim: {
    'jquery-ui': ['jquery'],
    'bootstrap': ['jquery-ui'],
    'd3': ['bootstrap'],
    'nv.d3': ['d3'],
    'stream': ['nv.d3'],
    'nv.d3.lib': ['stream'],
```

Current State

Mapa das Organizações da Sociedade Civil

Mapa das Organizações da Sociedade Civil

Organização

Município

Estado

Região

vasco

Buscar



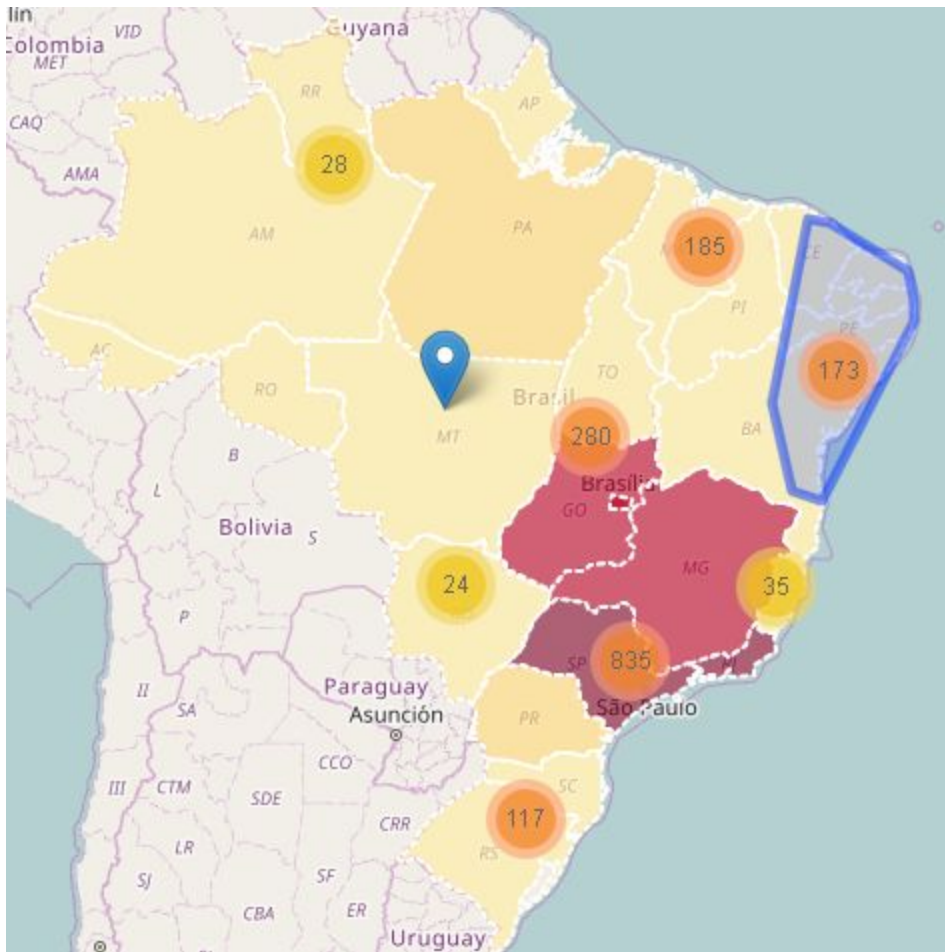
Map

Heatmap

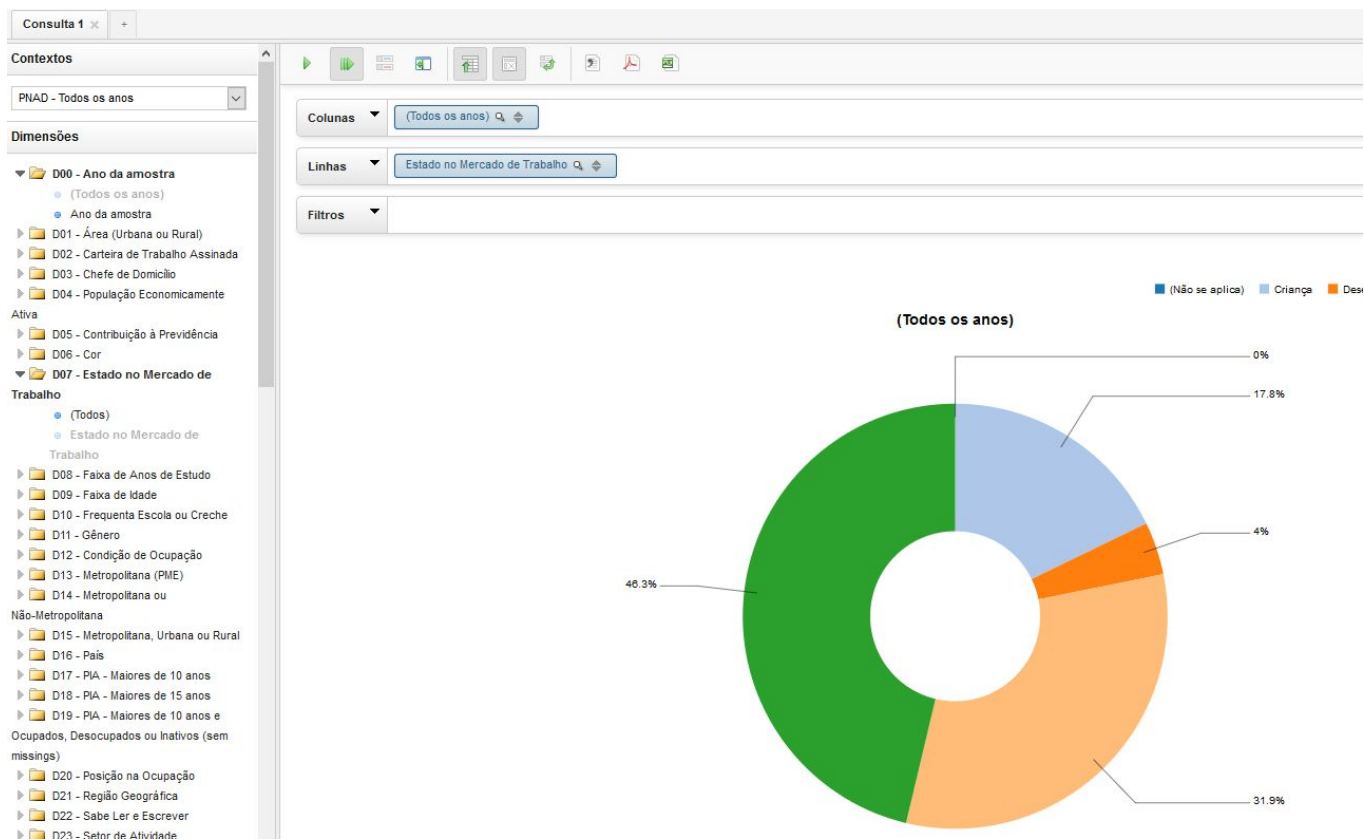
Clustering

Cache

Data sent by chunks



BI



2016/10/21

Result list

Resultado da Consulta

Visualizar por Mapa

Visualizar por Lista

Show 10 entries

Search:

NOME DA OSC



CNPJ



NATUREZA JURÍDICA



ENDEREÇO

DETALHAR

Vasco Da Gama
F C

87558110000151

Associação Privada

AV OSVALDO ARANHA,
1448, CIDADE ALTA,
Bento Gonçalves,
95700000

Detalhar



Gremio Dram E

Rectivo Vasco Da
Gama F C

46529723000183

Associação Privada

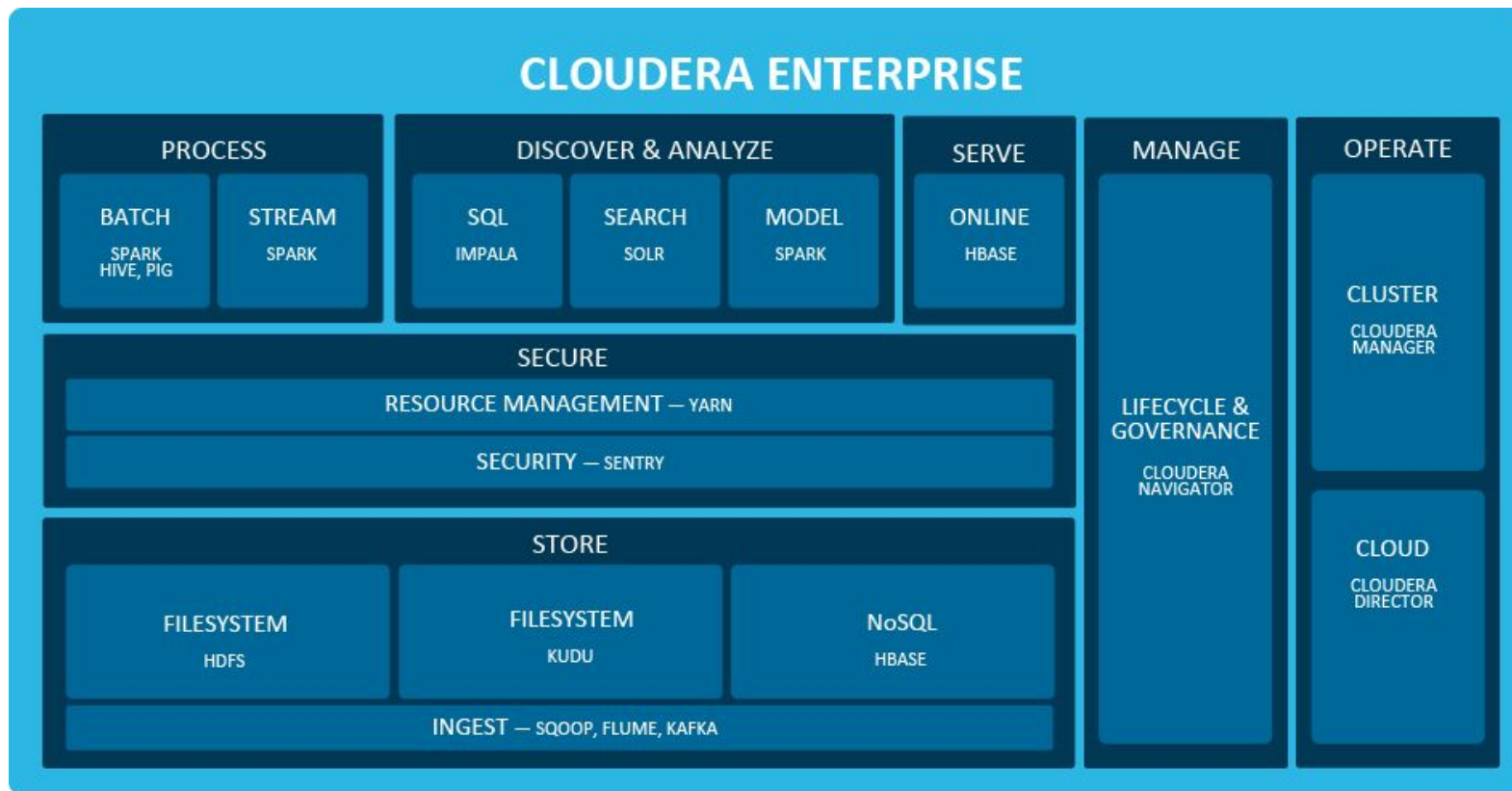
RUA JOSE BERNARDO
PINTO, 486, VILA
GUILHERME, São Paulo,
2055000

Detalhar

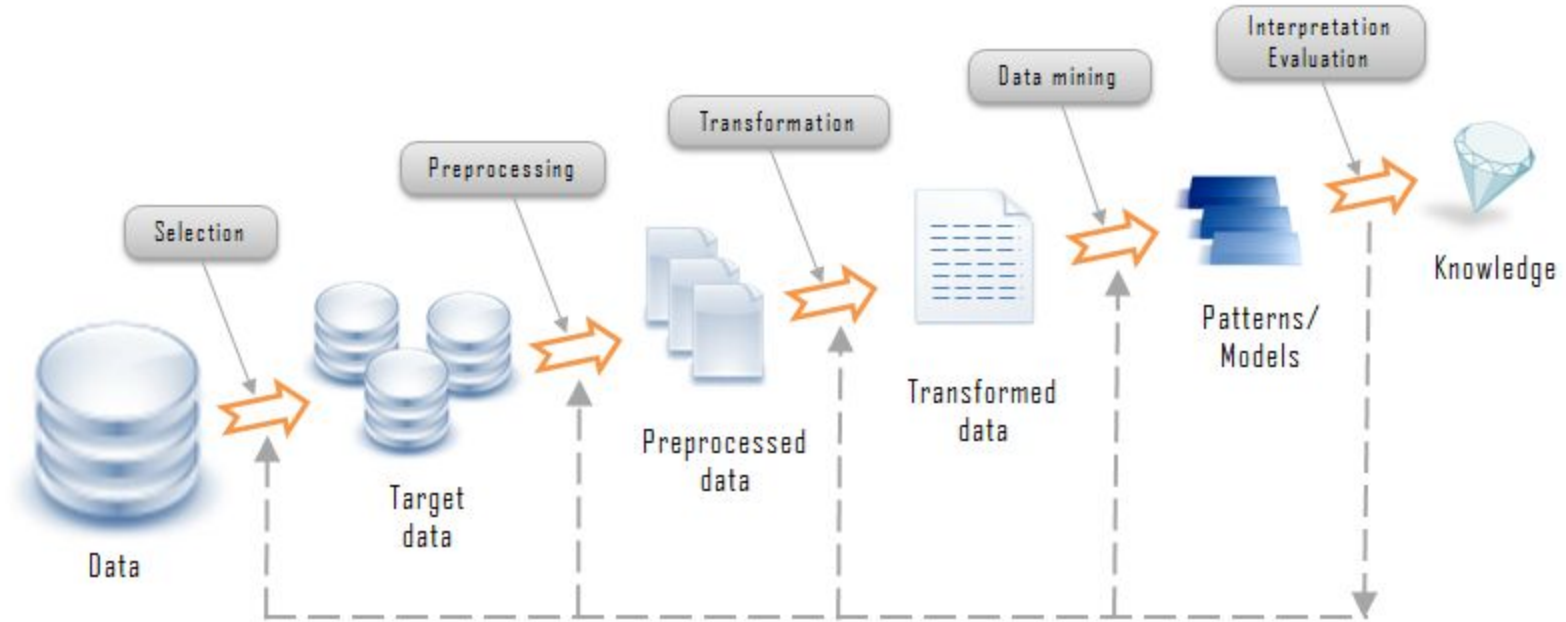


Future Steps

Big data full stack



Data mining & Knowledge discovery



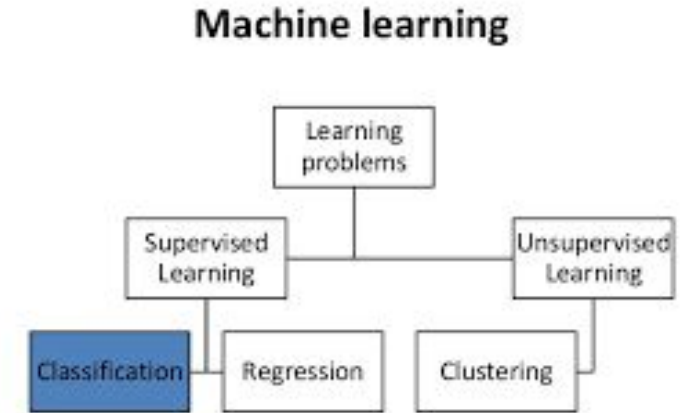
Applied machine learning & Information retrieval

Data classification (labeling)

Text searching and extraction

Time series prediction

Learning on graphs





Get in Touch

www.raulferreira.com.br

<https://br.linkedin.com/in/raulsenaferreira>

raul.ferreira@ipea.gov.br

raulsf@cos.ufrj.br

Open calls:

http://www.ipea.gov.br/portal/index.php/?option=com_content&view=article&id=23255&Itemid=5

A word cloud featuring the phrase "Thank You" in numerous languages and scripts. The words are arranged in a circular pattern, with "THANK YOU" being the largest and most central. Other languages include Spanish (GRACIAS, ARIGATO), Japanese (ARIGATO, SHUKURIA), Chinese (THANKS, 谢谢), and many others. The words are in different colors and orientations, creating a vibrant and multicultural display.