# Análise de Dados com Base em Processamento Massivo em Paralelo

# Twitter Data Lake

André Perez
ICMC/USP
andre.marcos.perez@usp.br
linkedin.com/in/andremarcosperez/

# Agenda

- Contexto

- Objetivo

- Arquitetura

- Links

# Agenda

- Contexto

- Objetivo

- Arquitetura

- Links

Política dos EUA no Twitter

# Congress Tweets Automator

```json
{
    "name":"Paul Cook",
    "chamber":"house",
    "type":"member",
    "party":"R",
    "accounts":[
        {
            "account_type":"campaign",
            "screen_name":"joinpaulcook",
            "id":"57177310"
        },
        {
            "id":"1074412920",
            "screen_name":"RepPaulCook",
            "account_type":"office"
        }
    ],
    "id":{
        "bioguide":"C001094",
        "govtrack":412513
    },
    "state":"CA"
}
```

COOK
COUNTY SUPERVISOR

Follow

Col. **@joinp**

Person
Superv

For offi

© App

142 Fo

Not follo

Follow

Rep. Paul Cook ✔
@RepPaulCook

Twenty-six year retired @USMC Colonel serving as a United States
Congressman from California's 8th Congressional District. Join the
conversation #CA08 #Semperfi

© Apple Valley, California  🔗 cook.house.gov  📅 Joined January 2013

**1,123** Following  **19K** Followers

Not followed by anyone you're following

# Congress Tweets

```
{
    "id":"1315057755430977539",
    "screen_name":"RepJacobs",
    "user_id":"1276232539510919168",
    "time":"2020-10-10T18:32:56-04:00",
    "link":"https://www.twitter.com/SPECNewsBuffalo/
    "text":"RT @SPECNewsBuffalo The Challenger Learr
    "source":"Twitter for iPhone"
}
```

t⫘ Rep. Chris Jacobs Retweeted

**Spectrum News BUF** ✔ @SPECNewsBuffalo · Oct 10
The Challenger Learning Center in Lockport uses space exploration as a theme to get kids interested in science, technology, engineering, and math.

Challenger Learning Center Wish Fulfilled With New Flag
The flag previously flown over the U.S. Capitol.
🔗 spectrumlocalnews.com

💬 1          t⫘ 1          ♡ 3

# Agenda

- Contexto

- **Objetivo**

- Arquitetura

- Links

Construir um pipeline de dados
que ingere dados diariamente
em um data lake
para armazenamento e exploração.

# Agenda

- Contexto

- Objetivo

- **Arquitetura**

- Links

# Arquitetura: Elementos

| ELT | Data Lake | Motor de Consulta | Orquestração |
|:---:|:---:|:---:|:---:|
|  |  |  |  |
| AWS Lambda | AWS s3 | AWS Athena | AWS Step Functions |
| | |  |  |
| | | AWS Glue | AWS CloudWatch |

# Arquitetura: Pipeline

Agenda

- Contexto

- Objetivo

- Arquitetura

- Links

# Links

Fonte de dados

- [Congress Tweets Automator](#)
- [Congress Tweets](#)

Amazon Web Services

- [AWS Service Docs](#)
- [AWS Simple Calculator](#)
- [AWS Newest Calculator](#)
- [AWS to Azure (Microsoft) Map](#)
- [AWS to GCP (Google) Map](#)

# Links

## Repositório de códigos

- [Data Science MBA @ ICMC-USP](#)

## Artigos

- [Demystify Hadoop Data Formats: Avro, ORC, and Parquet](#)
- [New in Hadoop: You should know the Various File Format in Hadoop](#)