# THE DEVELOPER'S CONFERENCE

# Trilha – BigData

## DataOps: Estendendo as práticas de DevOps para BigData

**Eduardo Hahn**

**DataLakers founder & DataOps Enthusiastic**

# + BuzzWord

> "**DataOps** is an automated, process-oriented methodology, used by analytic and data teams, to improve the quality and reduce the cycle time of data analytics." Wikipedia

> "**DataOps** is about more than speed and quality. With a culture of continuous improvement, organizations can deliver data analytics solutions more efficiently, releasing valuable team members for more valuable activities, such as building innovative new products." Eckerson Group

# Começo.....

> 2014....

**Big Data & Analytics Hub** (IBM)

## Blogs

### 3 reasons why DataOps is essential for big data success

JUNE 19, 2014

by Lenny Liebmann
Contributing Editor, InformationWeek
Follow me on LinkedIn, Twitter

Developers once wrote application code and just "threw it over the wall" to IT operations, which then had to ensure that those applications performed well in the production environment. This was always a less-than-optimal approach, but it became untenable as the business began to depend more and more on lots of fresh code getting rolled out into production quickly and with a high degree confidence. So IT organizations are now embracing a set of best practices known as *DevOps* that improve coordination between development and operations.

"A ciência de dados é uma disciplina excepcionalmente importante hoje em dia. Mas essa ciência só é útil na medida em que pode ser executada de forma eficiente e confiável. E para que isso aconteça, você precisa de DataOps."

DataLakers
The Big Data Company

# Começo.....

> 2015...



⁘ tamr    Products ⌄    Services    Solutions ⌄    Customers    Resources ⌄    Company ⌄    Contact

# From DevOps to DataOps, By Andy Palmer

Posted on Thursday, May 7th, 2015 at 1:55 PM.

Written by Andy Palmer

## Why It's Time to Embrace "DataOps" as a New Discipline

Over the past 10 years, many of us in technology companies have experienced the emergence of "DevOps." This new set of practices and tools has improved the velocity, quality, predictability and scale of software engineering and deployment. Starting at the large internet companies, the trend towards DevOps

**DataLakers**
The Big Data Company

# Acelerou.....

> 2017....

# Agora sim.....

> 2018...

Figure 1. Hype Cycle for Data Management, 2018

# Agora sim.....

DataOps is a collaborative data management practice focused on improving the communication, integration and automation of data flows between data managers and data consumers across an organization.

Gartner.

# DataOps Manifesto

## The DataOps Manifesto

Through firsthand experience working with data across organizations, tools, and industries we have uncovered a better way to develop and deliver analytics that we call DataOps.

Whether referred to as data science, data engineering, data management, big data, business intelligence, or the like, through our work we have come to value in analytics:

Individuals and interactions over processes and tools
Working analytics over comprehensive documentation
Customer collaboration over contract negotiation
Experimentation, iteration, and feedback over extensive upfront design
Cross-functional ownership of operations over siloed responsibilities
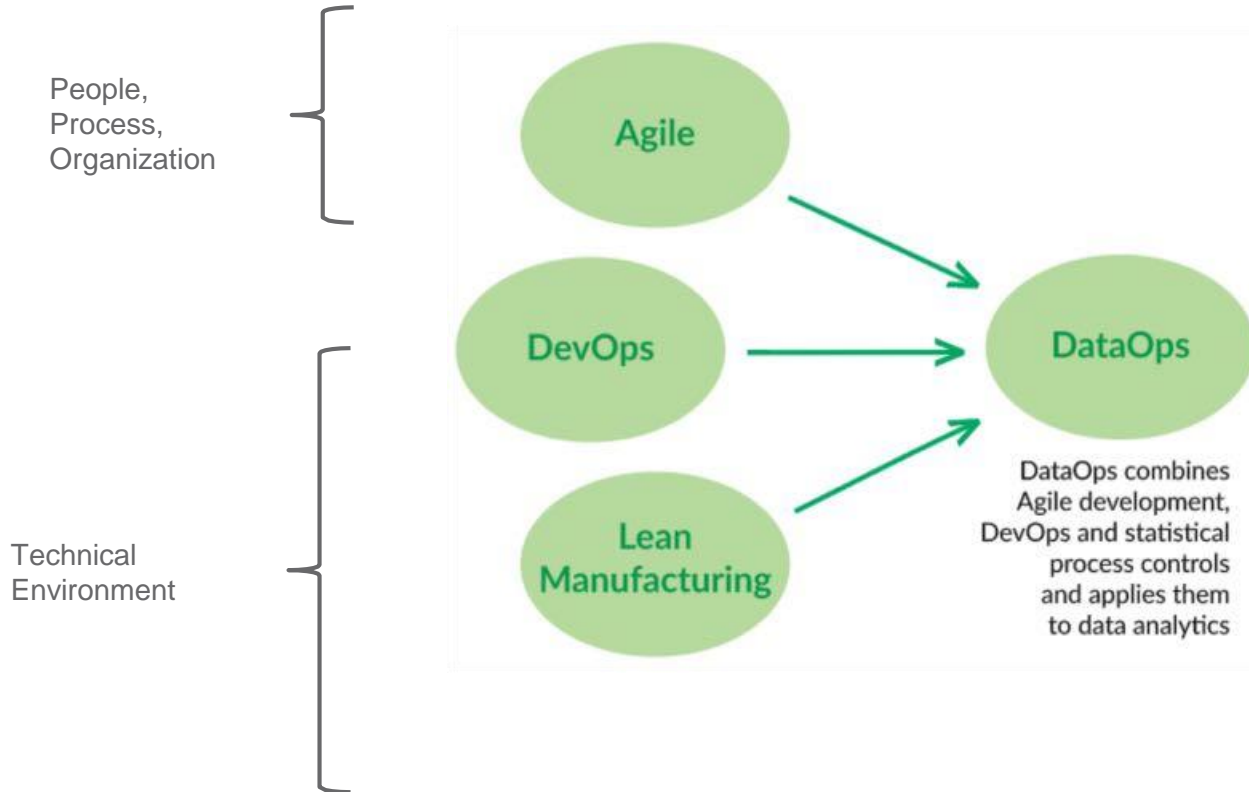
dataopsmanifesto.org/

# DataOps Manifesto

**Princípios do DataOps**

- **1. Satisfaça continuamente o seu cliente**
- **2. Valor do trabalho analítico**
- **3. Abrace a mudança**
- **4. É um esporte em equipe**
- **5. Interações diárias**
- **6. Auto-organização**
- **7. Reduza o heroísmo**
- **8. Reflita**
- **9. Os códigos**

- **10. Orquestração**
- **11. Faça tudo ser reproduzível**
- **12. Ambientes descartáveis**
- **13. Simplicidade**
- **14. Análise de dados é manufatura**
- **15. A qualidade é primordial**
- **16. Monitorar a qualidade e o desempenho**
- **17. Reutilizar**
- **18. Melhorar os tempos dos ciclos**

# Genesis of DataOps



People, Process, Organization

Technical Environment

Agile

DevOps

Lean Manufacturing

DataOps

DataOps combines Agile development, DevOps and statistical process controls and applies them to data analytics

# 4 "As" de DataOps

> Automatize e monitorar pipelines

> Automatizar implantações

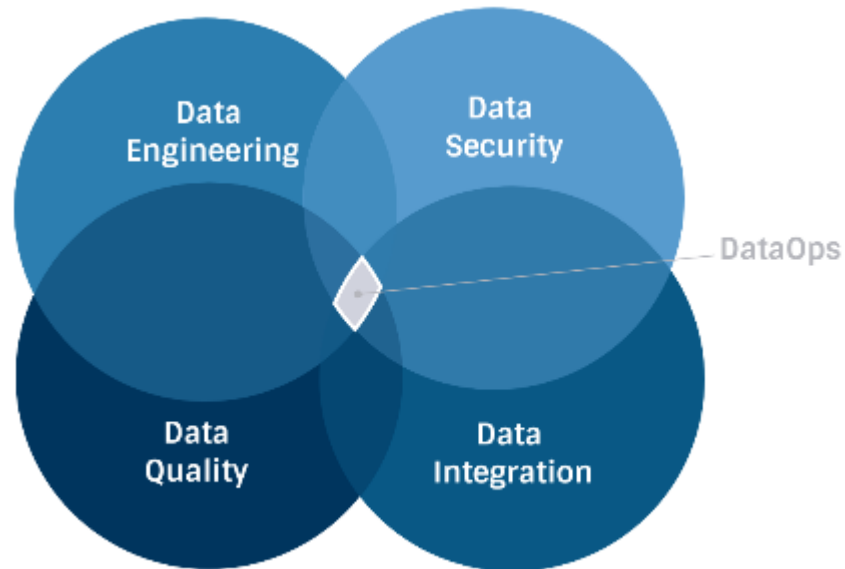> Automatizar e monitorar a qualidade

> Automatizar sandbox
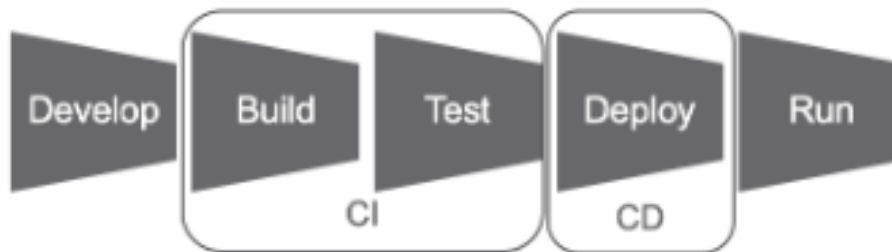
# DataOps is NOT Just DevOps for Data



**DEVOPS**

**DATAOPS**

# DataOps is NOT Just DevOps for Data

# DataOps is NOT Just DevOps for Data

# DataOps is NOT Just DevOps for Data

**Automatize e monitorar pipelines**

# DataOps is NOT Just DevOps for Data

**Automatizar implantações**

# DataOps is NOT Just DevOps for Data

## Automatizar sandbox

# DataOps is NOT Just DevOps for Data

**Automatizar e monitorar a qualidade**

# The People of DataOps

# The People of DataOps

# The People of DataOps

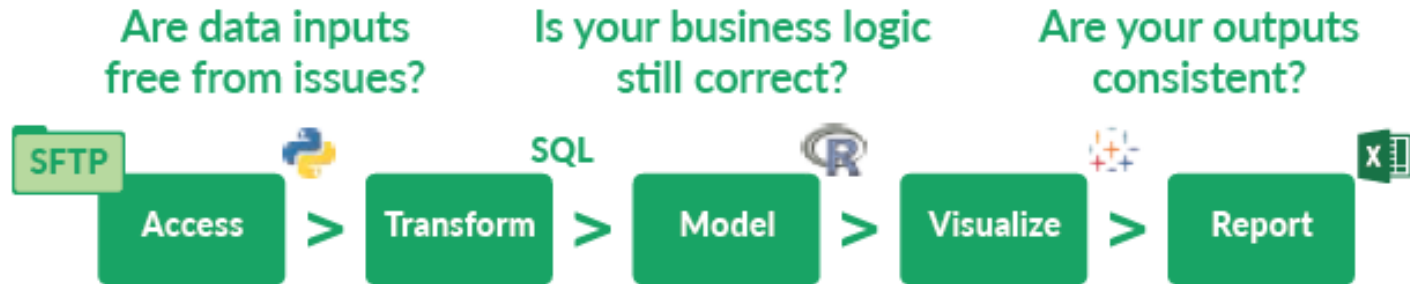| | Role | Goals | Tools |
|---|---|---|---|
| **Consumers** | **Citizen** | Use data to make business decisions | Viz, CRM, Excel, PowerPoint, Word, Web Search |
| | **Analyst** | Deliver insights to the business, typically through dashboards and reports | Viz, Excel, SSDP, Web Search |
| | **Scientist** | Deliver insights to the business, typically through models and algorithms | R, Python, SAS, SSDP |
| | **Developer** | Build applications which leverage corporate data | Python, Java, JS, SQL, REST |
| **Preparers** | **Engineer** | Deliver and manage data pipelines | ETL, SQL, Python |
| | **Curator** | Ensure consumers have the data they need, in the form they need it | Mastering tools, Catalog |
| | **Steward** | Use feedback from consumers to improve data broadly, ensure governance | Feedback tools, Governance |
| **Suppliers** | **Source Owner** | Define and manage purpose, processes (data creation, consumption) & users (i.e., access) of the data source | EDW, SQL, ERWin, LDAP, SAP |

# Passos para implementar DataOps

# Passos para implementar DataOps
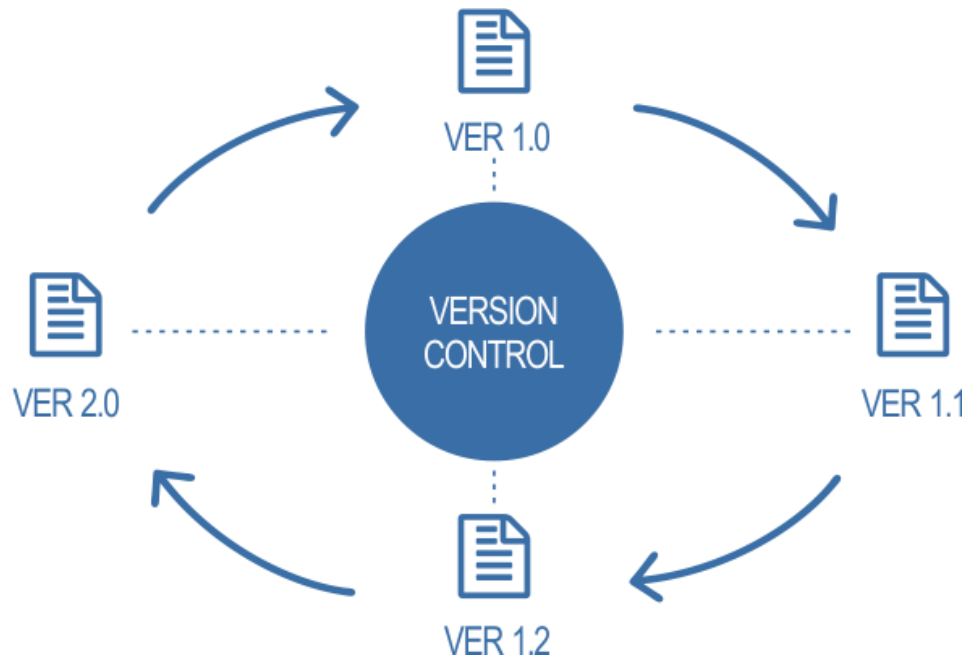
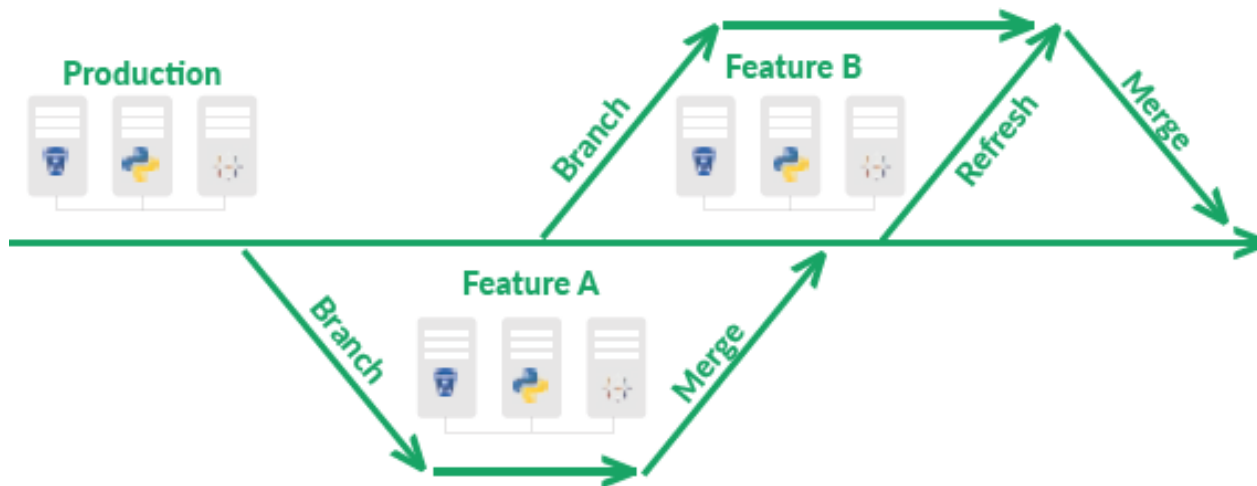> Passo 1: implemente Data Test e Logic Test

# Passos para implementar DataOps

> Passo 2: use controle de versão (Branch&Merge).

# Passos para implementar DataOps

> ## Passo 3: Use múltiplos ambientes

# ❯ Passo 4: Reuso e containers

Outro método de aumento de produtividade para equipes é a capacidade de reutilizar e conter o código. Cada etapa intermediária no pipeline de análise de dados recebe a saída de um estágio anterior e fornece entrada para o próximo estágio. É mais fácil para os outros membros da equipe reutilizar componentes menores, se eles puderem ser segmentados ou conteinerização. Melhor caminho é usar Docker

## Passo 5: Parametrize seus processos

# Passo 6: Use Simple Unique Storage

- Data Lake: os dados são movidos de diferentes silos de dados para um repositório comum, é muito mais fácil para uma equipe de análise de dados trabalhar com ele.
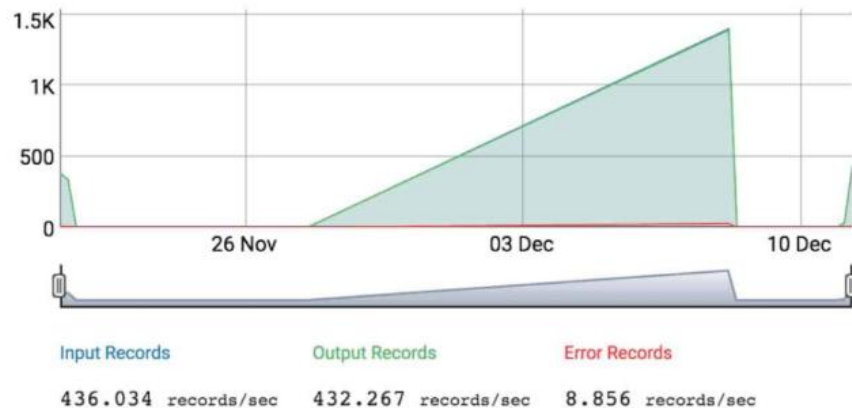
## Passo 7: Medir todas as etapas do pipeline

Defina métricas de ponta a ponta para sua arquitetura. Identifique pontos de melhorias e problemas de desempenho. Visualize uma arquitetura de dados online, visualizando como os sistemas evoluem.
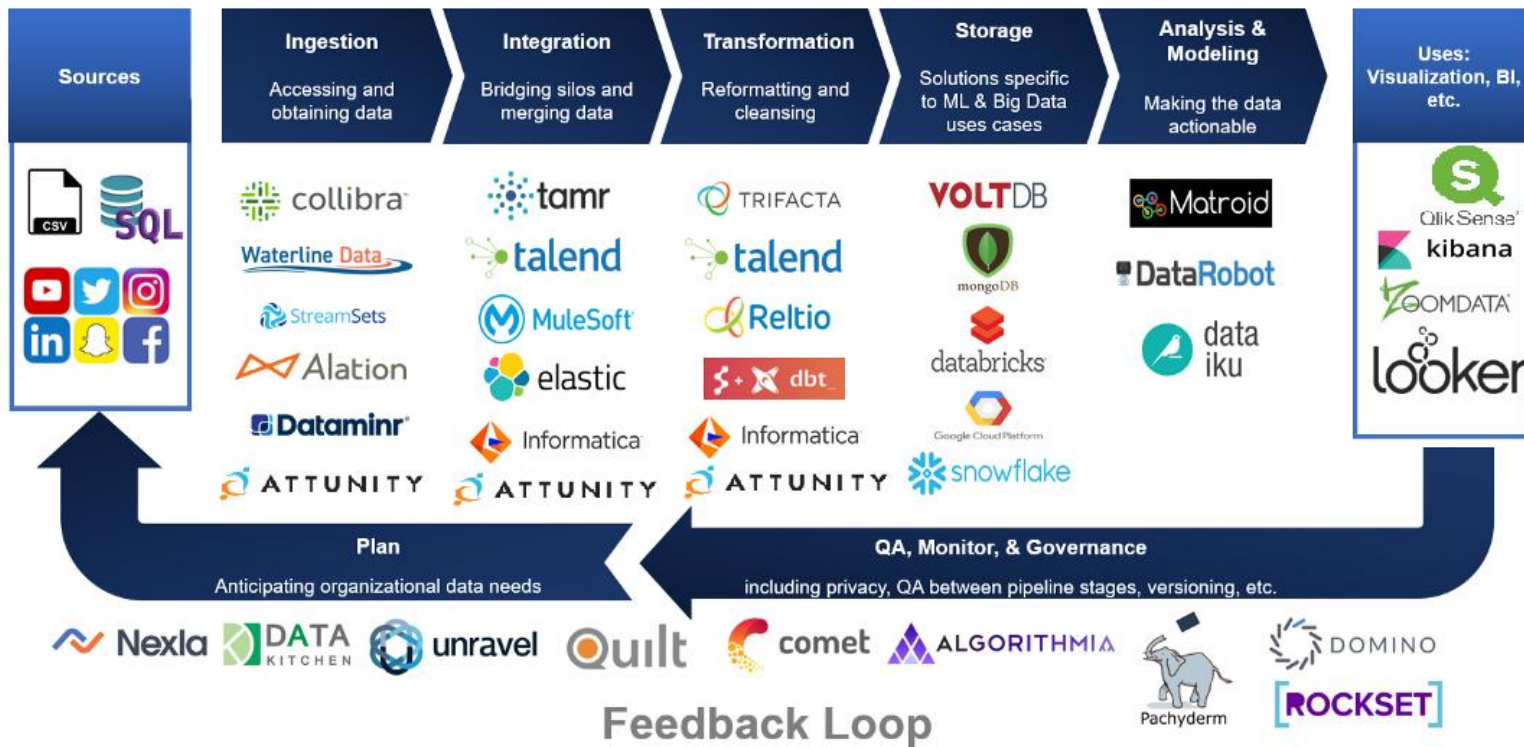


Record Throughput Time Series

| Input Records | Output Records | Error Records |
|---|---|---|
| 436.034 records/sec | 432.267 records/sec | 8.856 records/sec |

# DataOps Ecosystem

› **Platform Solutions**

# DataOps Ecosystem



Data Delivery Pipeline, from Sources to Uses

# References



**O'REILLY®**

Creating a
Data-Driven Enterprise
with DataOps

Insights from Facebook, Uber, LinkedIn,
Twitter, and eBay

Ashish Thusoo &
Joydeep Sen Sarma



THE
DEVELOPER'S
CONFERENCE



**DataLakers**
The Big Data Company

# References

- **DataOps Ecosystem - //**medium.com/data-ops/2017-the-year-of-dataops-b2023c17d2af
- **DataOps for Government (State of Connecticut) -** //github.com/OpenDataCT/DataOps
- **Creating a Data-Driven Enterprise with DataOps - //**www.oreilly.com/data/free/creating-a-data-driven-enterprise-with-dataops.csp
- http://dataopsmanifesto.org/
- **DataOps—It's a Secret** - //www.datasciencecentral.com/profiles/blogs/dataops-it-s-a-secret
- **The Power of DataOps** - //www.delphix.com/blog/power-dataops
- **Building a DataOps Team** - //medium.com/data-ops/building-a-dataops-team-abc375e0a6bc
- **DataOps: Industrializing Data and Analytics -** //www.eckerson.com/articles/dataops-industrializing-data-and-analytics?content=dataops-industrializing-data-and-analytics

# Concluindo....

> As empresas que desejam implementar DataOps devem concentrar seus esforços em três áreas:
>> Cultura
>> Processos
>> Tecnologia

# Oportunidades estão chegando....
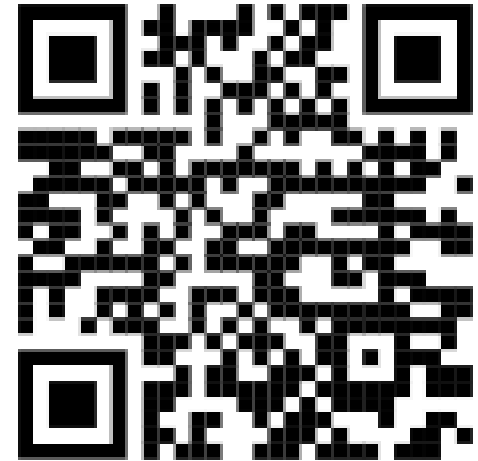
# About me

❯ Eduardo Hahn

- Founder DataLakers Tecnologia
- Data Lover & DataOps Enthusiastic
- eduardo.hahn@datalakers.com.br
- @eduardohahn
- /in/eduardohahn3

Big Family

Partners

cloudera®  talend  ATTUNITY

DataLakers
The Big Data Company