



# THE DEVELOPER'S CONFERENCE

## Trilha – Big Data

**Passo a Passo para implementar DataOps em projetos de Big Data**

**Eduardo Hahn**



- + 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.....

IBM Big Data & Analytics Hub

## 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 of 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.”

# Começo.....



THE  
DEVELOPER'S  
CONFERENCE

> 2015...



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....



M  [Follow](#)  
HOME STORY INDEX ON QUALITY NEWS CDO 7 STEPS SPEED ML | [LEARN MORE AT DATAKITCHEN.IO](#)

 DataKitchen  
May 9, 2017 · 4 min read

### The DataOps Ecosystem Emerges

In 2015, Andy Palmer of Tamr defined the term **DataOps**, a faster, more flexible approach to data analytics which recognizes the interconnectedness of IT operations, data engineering, data integration, data quality and data security/privacy.



**FEATURE**

## What is DataOps? Collaborative, cross-functional analytics

DataOps (data operations) is an emerging discipline that brings together DevOps teams with data engineer and data scientist roles to provide the tools, processes and organizational structures to support the data-focused enterprise.



By **Thor Olavsrud**  
Senior Writer, CIO | NOV 21, 2017 3:00 AM PT



**INSIDER** [Sign In](#) | [Register](#)

**BrandPost** Sponsored by Delphix | [Learn More](#)



## UNLOCKING INNOVATION WITH DATA

By **Eric Schrock** | NOV 9, 2017 6:24 AM PT

**SPONSORED**

### The Power of DataOps

The digital economy has created an unquenchable thirst for data across all aspects of business. Those that can leverage data to drive innovation will win.





Billionaires Innovation Lead



3,992 views | Apr 11, 2018, 01:38pm

## How DataOps Is Transforming Data Management Practices

 **Randy Bean** Contributor  
**CIO Network** Contributor Group

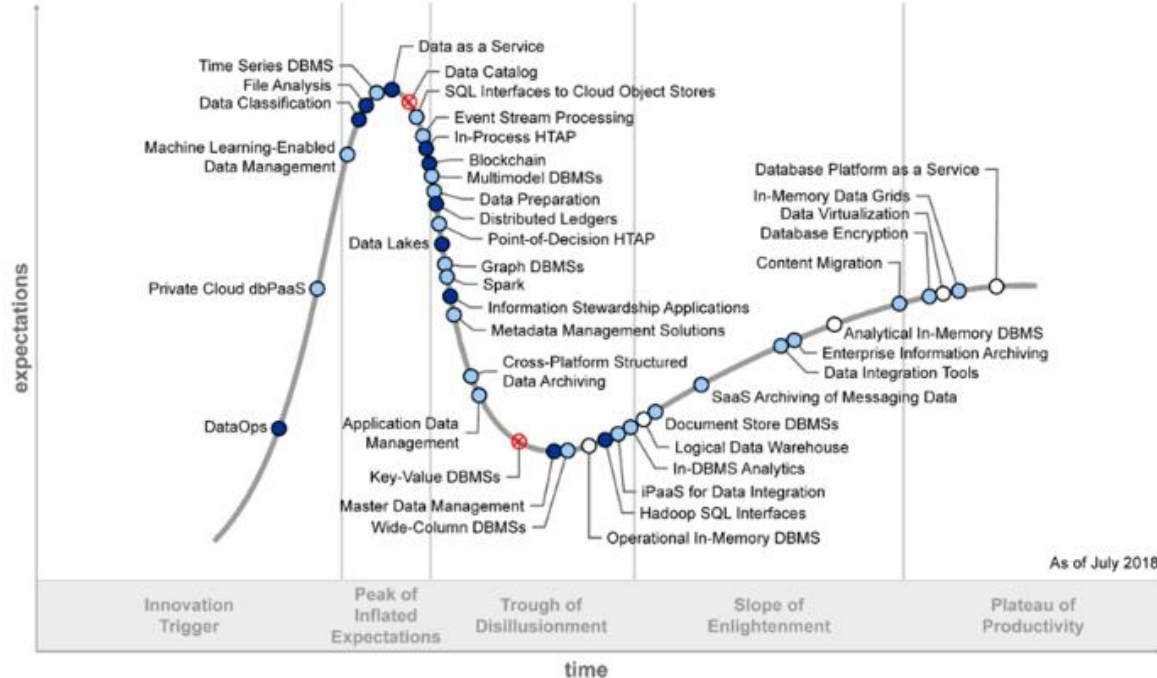
# Agora sim.....



THE  
DEVELOPER'S  
CONFERENCE

Figure 1. Hype Cycle for Data Management, 2018

➤ 2018...



Plateau will be reached:

- less than 2 years
- 2 to 5 years
- 5 to 10 years
- ▲ more than 10 years
- ⊗ obsolete before plateau

© 2018 Gartner, Inc.

# Agora sim.....



THE  
DEVELOPER'S  
CONFERENCE

**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  
DEVELOPER'S  
CONFERENCE

## 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/](https://dataopsmanifesto.org/)

**DataLakers**  
The Big Data Company



# DataOps Manifesto



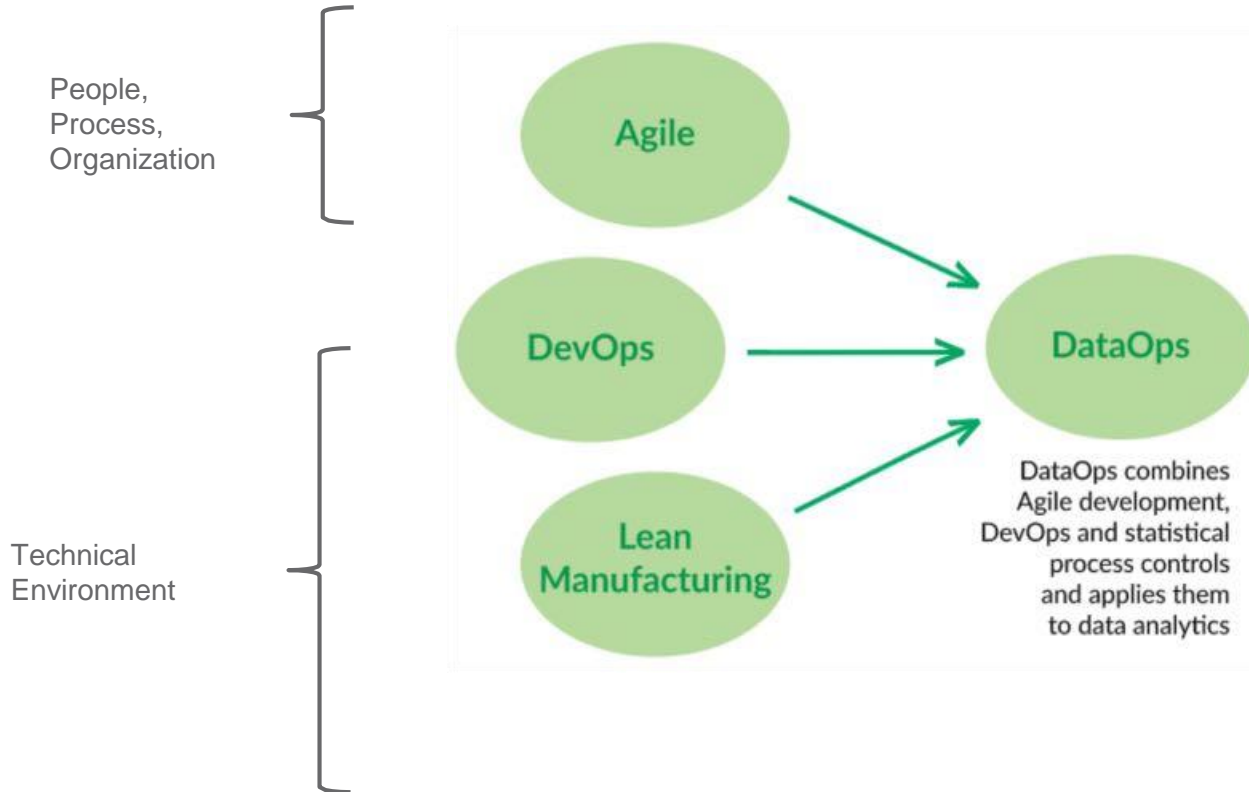
## ➤ Princípios do DataOps

- 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



THE  
DEVELOPER'S  
CONFERENCE



# 4 “As” de DataOps



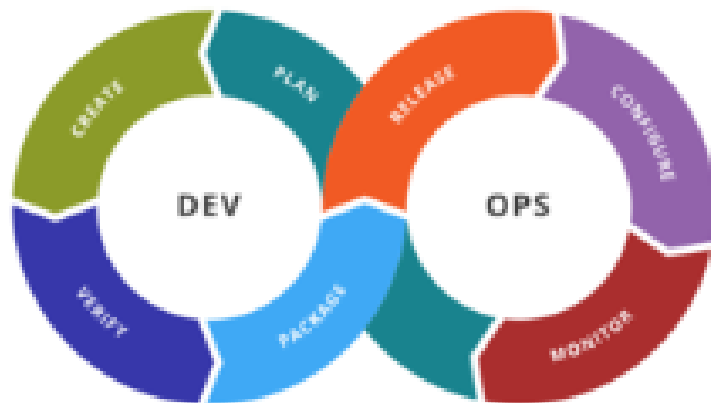
THE  
DEVELOPER'S  
CONFERENCE

- Automatize e monitorar pipelines
- Automatizar implantações
- Automatizar e monitorar a qualidade
- Automatizar sandbox

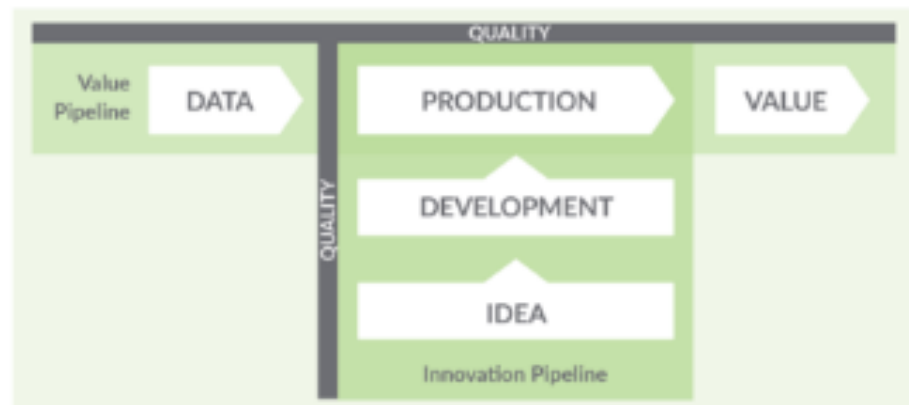
# DataOps is NOT Just DevOps for Data



THE  
DEVELOPER'S  
CONFERENCE



DevOps



DataOps

# DataOps is NOT Just DevOps for Data



THE  
DEVELOPER'S  
CONFERENCE

## Different People and Expectations

DevOps  
Users &  
Tools



Software Engineers, comfortable with coding and complexity of multiple languages, tools, and hardware/software.

DataOps  
User &  
Tools



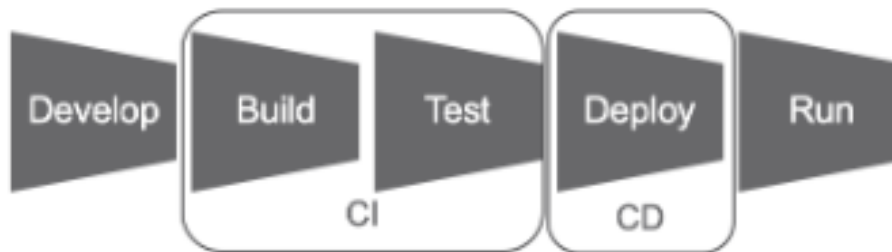
Data Scientists, Engineers, and Analysts who want to just analyze data and build models.

# DataOps is NOT Just DevOps for Data

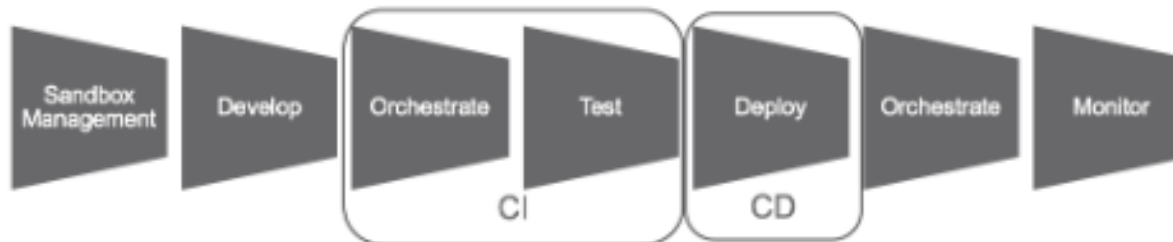


THE  
DEVELOPER'S  
CONFERENCE

**DevOps  
Process**



**DataOps  
Process**

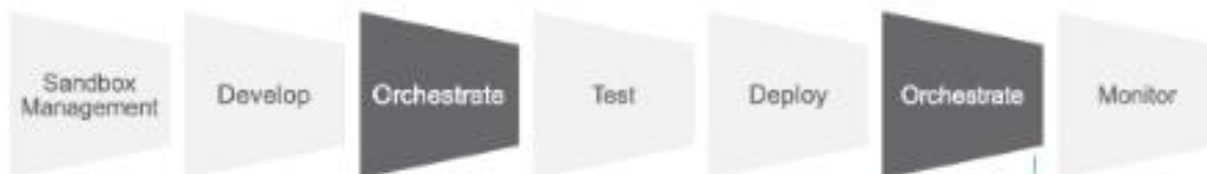


# DataOps is NOT Just DevOps for Data



THE  
DEVELOPER'S  
CONFERENCE

## DataOps Process



DataOps Production  
Requires Automated  
Testing, Monitoring and  
Statistical Process Control

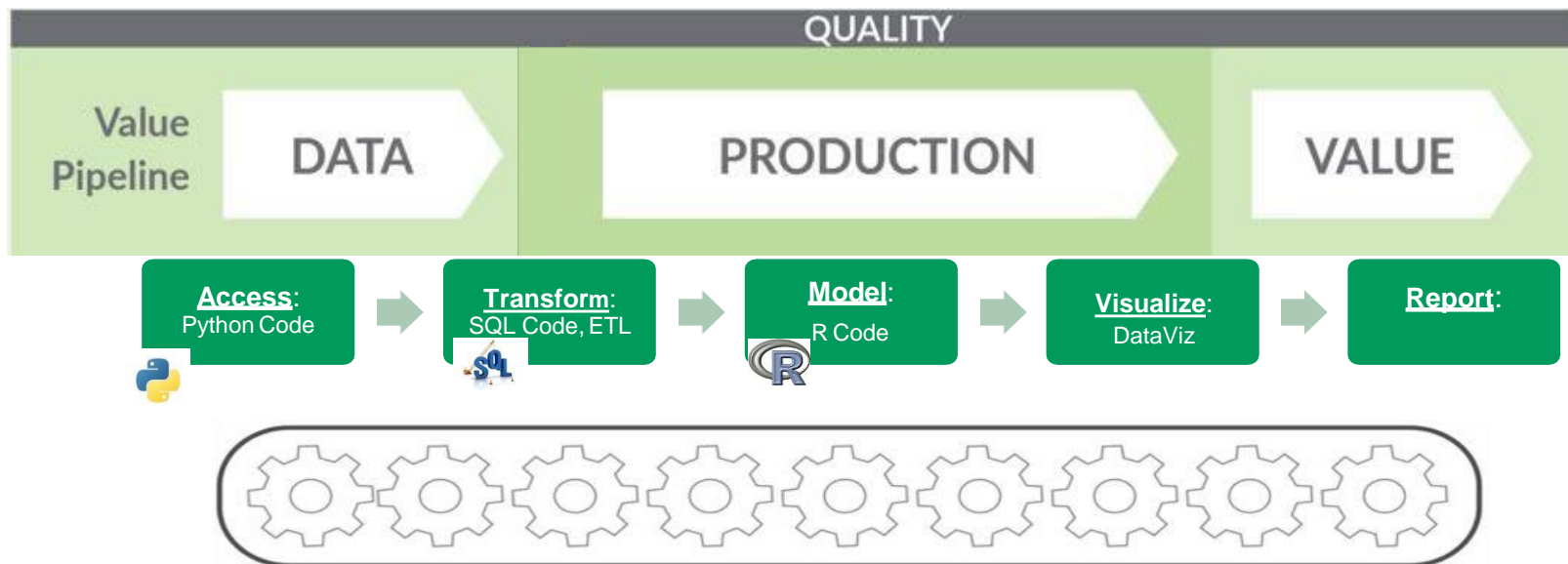


# DataOps is NOT Just DevOps for Data



THE  
DEVELOPER'S  
CONFERENCE

Automatize e monitorar pipelines



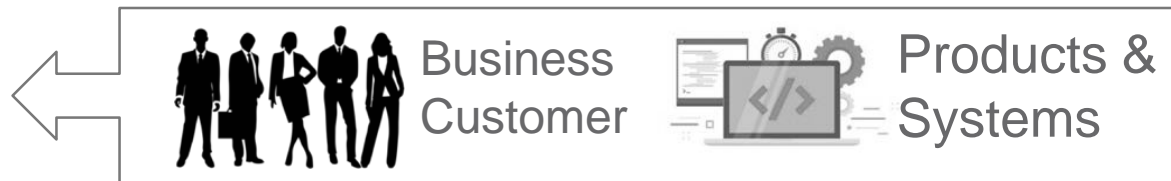
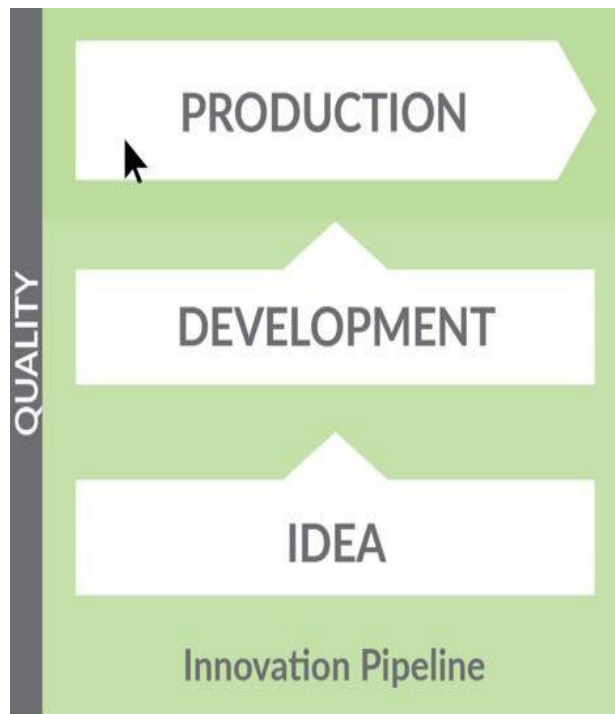


# DataOps is NOT Just DevOps for Data

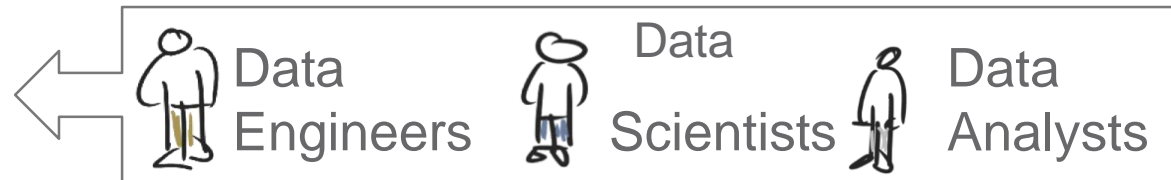


THE  
DEVELOPER'S  
CONFERENCE

## Automatizar implantações



**Diverse Team**



**Diverse Tools**



# DataOps is NOT Just DevOps for Data

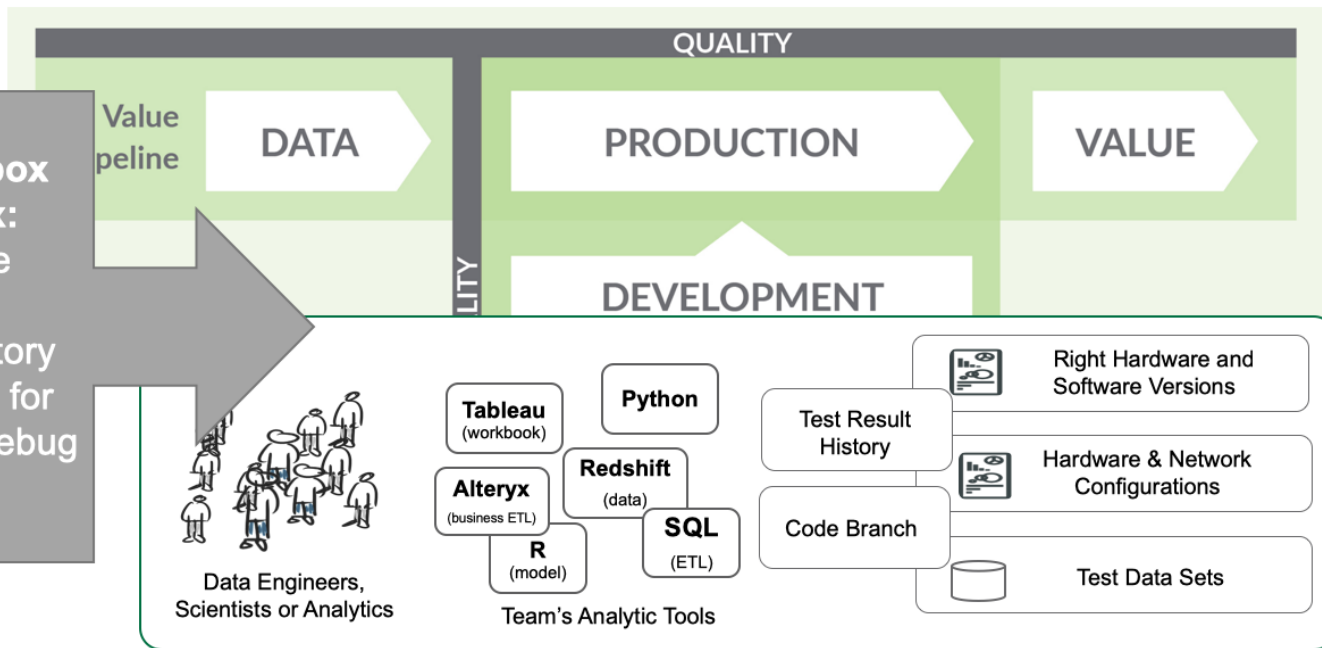


THE  
DEVELOPER'S  
CONFERENCE

## Automatizar sandbox

### Development Sandbox Creation is Complex:

Hard to create the right set of data, tools, people, history and configuration for a fast build test debug cycle

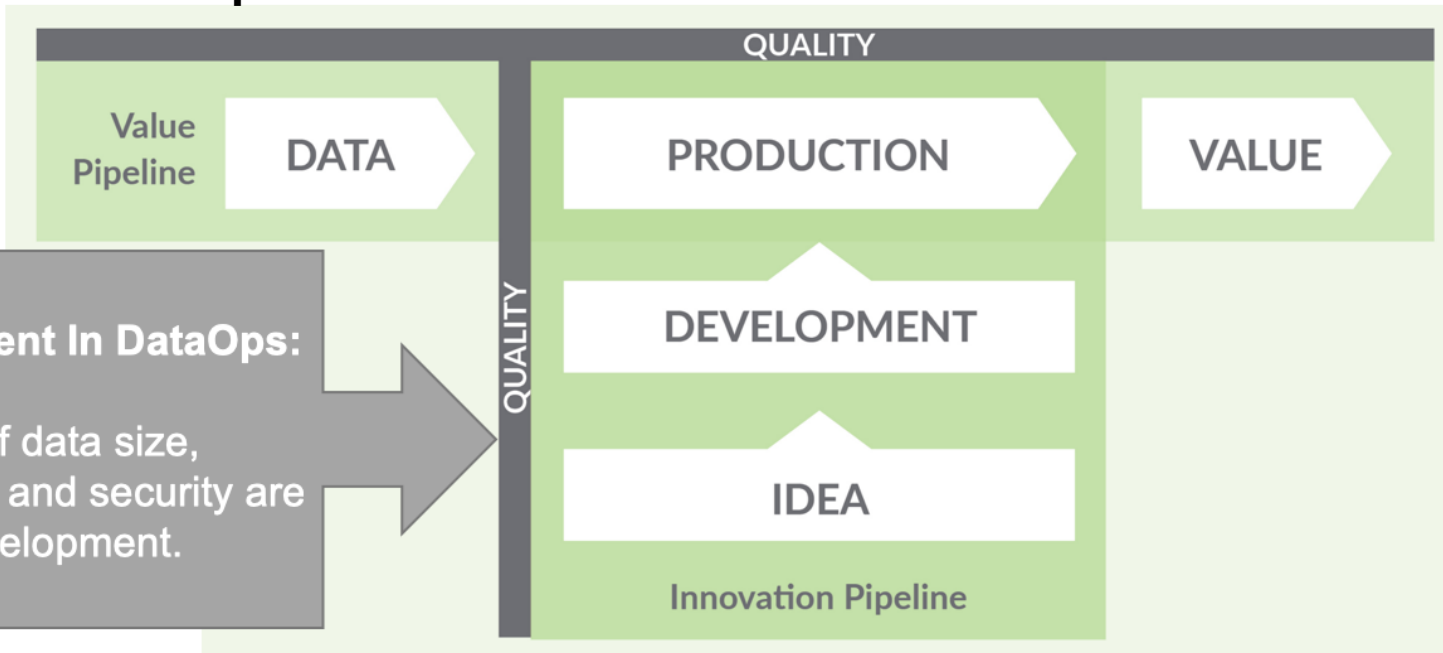


# DataOps is NOT Just DevOps for Data



THE  
DEVELOPER'S  
CONFERENCE

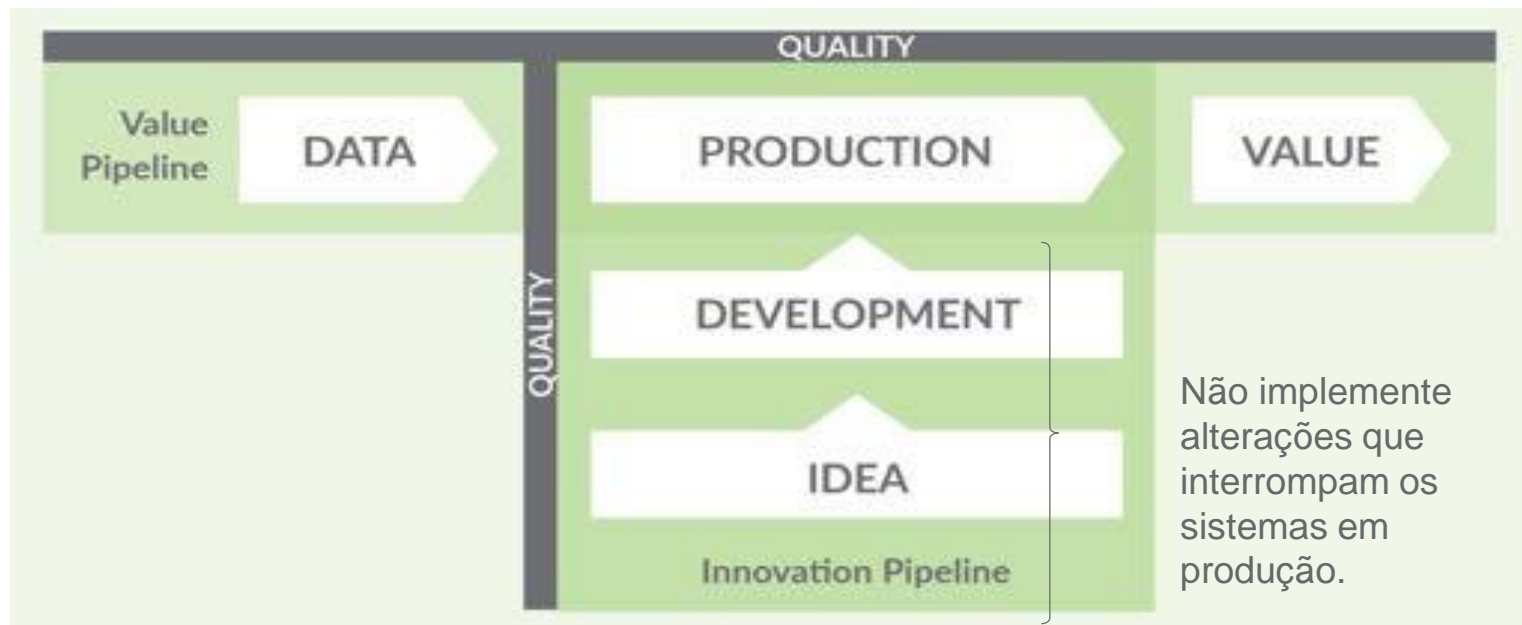
Automatizar e monitorar a qualidade



# DataOps is NOT Just DevOps for Data



Não permita que dados de baixa qualidade cheguem aos usuários no **Value Pipeline**



# Passos para implementar DataOps



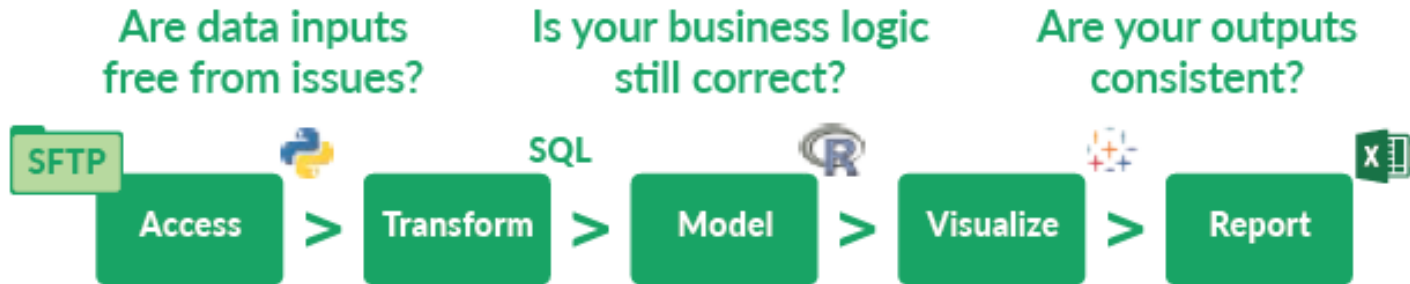
THE  
DEVELOPER'S  
CONFERENCE



# Passos para implementar DataOps



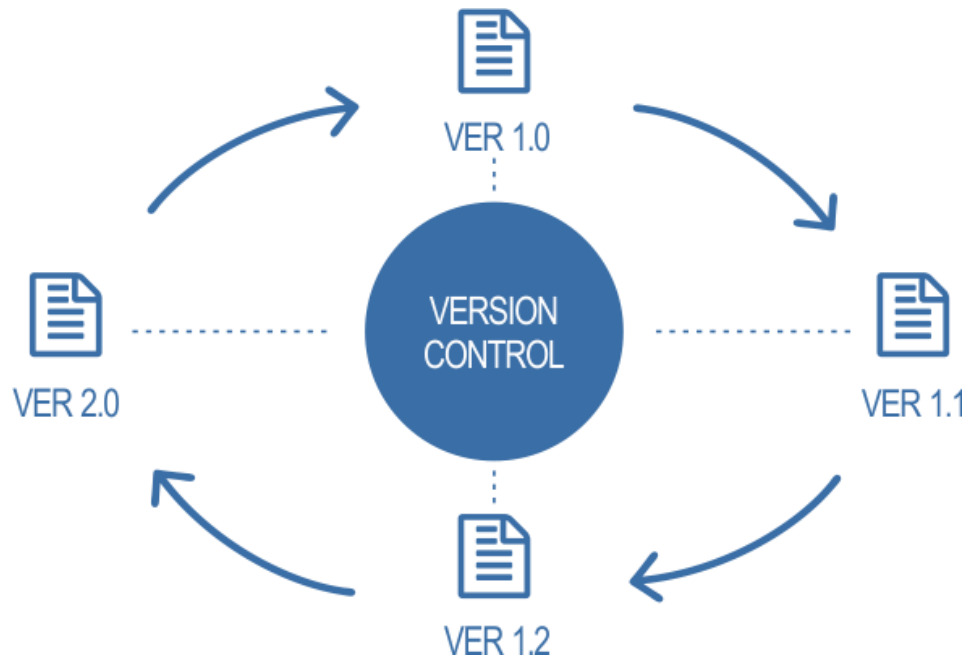
## ➤ Passo 1: incremente Data Test e Logic Test



# Passos para implementar DataOps

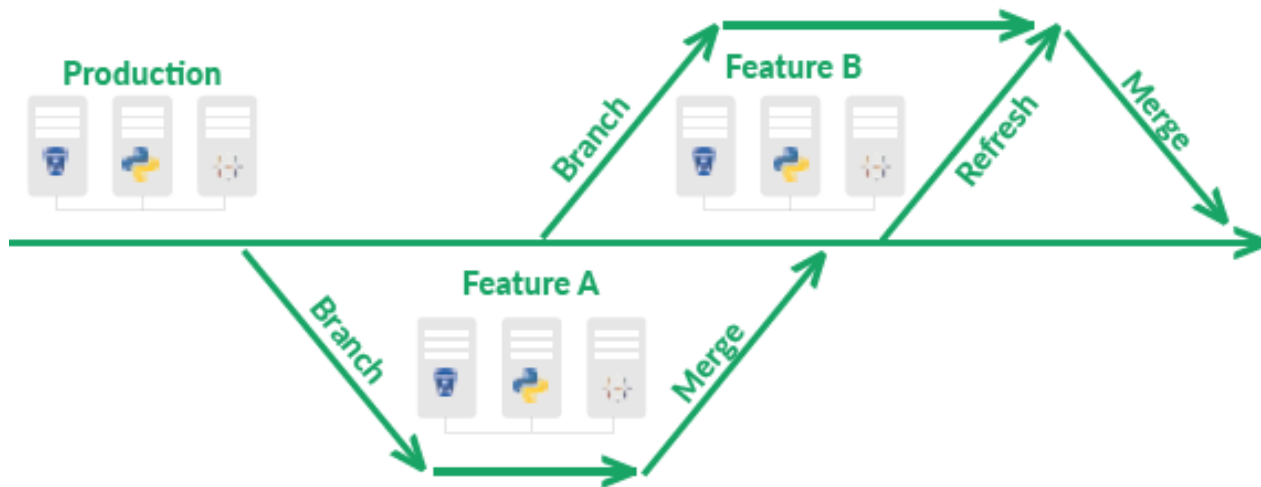


- Passo 2: use controle de versão (Branch&Merge).



# Passos para implementar DataOps

## ➤ Passo 3: Use múltiplos ambientes





# Passos para implementar DataOps



## ➤ 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 containerização. Melhor caminho é usar Docker



# Passos para implementar DataOps



THE  
DEVELOPER'S  
CONFERENCE

## ➤ Passo 5: Parametrize seus processos

This project is parameterized

**Git Parameter**

Name

Description

Parameter Type

# Passos para implementar DataOps



## ➤ 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.

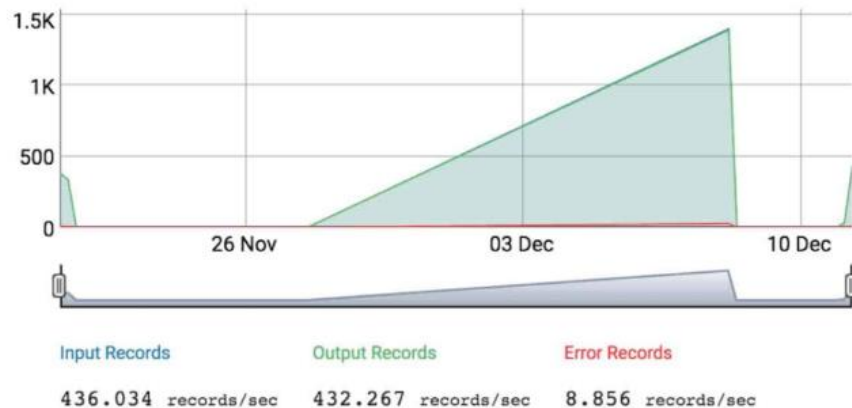
# Passos para implementar DataOps



## ➤ 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



# DataOps Ecosystem



## > Platform Solutions

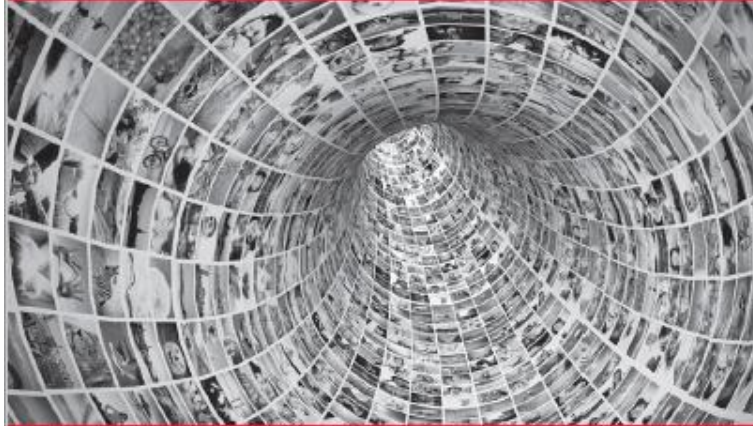


# References

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](https://medium.com/data-ops/2017-the-year-of-dataops-b2023c17d2af)
- **DataOps for Government (State of Connecticut)** - [//github.com/OpenDataCT/DataOps](https://github.com/OpenDataCT/DataOps)
- **Creating a Data-Driven Enterprise with DataOps** - [//www.oreilly.com/data/free/creating-a-data-driven-enterprise-with-dataops.csp](https://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](https://www.datasciencecentral.com/profiles/blogs/dataops-it-s-a-secret)
- **The Power of DataOps** - [//www.delphix.com/blog/power-dataops](https://www.delphix.com/blog/power-dataops)
- **Building a DataOps Team** - [//medium.com/data-ops/building-a-dataops-team-abc375e0a6bc](https://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](https://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



# About me



Big Family

- Eduardo Hahn
  - Founder DataLakers Tecnologia
  - Data Lover & DataOps Enthusiastic
  - [eduardo.hahn@datalakers.com.br](mailto:eduardo.hahn@datalakers.com.br)
  - [@eduardohahn](https://twitter.com/eduardohahn)
  - [/in/eduardohahn3](https://www.linkedin.com/in/eduardohahn3)



Partners



DataLakers  
The Big Data Company



# THE DEVELOPER'S CONFERENCE

## Trilha – Big Data

**Passo a Passo para implementar DataOps em projetos de Big Data**

**Eduardo Hahn**

**DataLakers founder & DataOps Enthusiastic**

**DataLakers**  
The Big Data Company