

# ANÁLISE DE DADOS UTILIZANDO CLUSTER DE BAIXO CUSTO

## COMPARAÇÃO DE DESEMPENHO DE AMBIENTES VIRTUAIS

Felipe Fonseca Rocha

Orientador: Ítalo Fernando Scotá Cunha

Universidade Federal de Minas Gerais

09 de Fevereiro de 2022



# Sumário





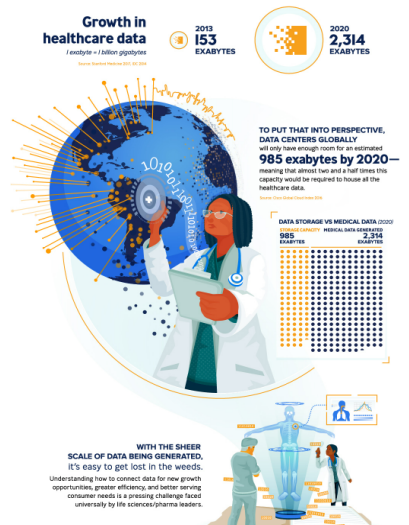
# Contexto e Motivação I



A todo momento nós geramos milhões de dados que são coletados por diferentes meios

Várias ferramentas estão disponíveis para Transformá-los em informações e embasar decisões

# Contexto e Motivação II



Isso também acontece na área da saúde  
Porém o uso dessas ferramentas nessa  
área, para transformar dados em in-  
formação, ainda é pouco significativo

## Contexto e Motivação III

- : Tendência crescente de trabalho interdisciplinar
- : Potencial de melhora do sistema de Saúde através de análise de dados
- : Necessário Propor e Validar estratégias que sejam viáveis e facilitem o processamento de análise de grande volume de dados produzido na área



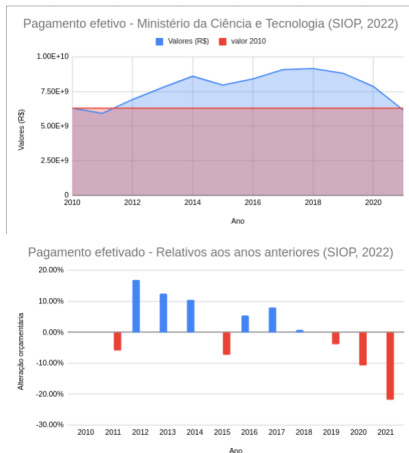
## Contexto e Motivação IV

- : No Brasil, dados do Sistema de Informação em Saúde (SIS) são disponibilizados desde 2016
- : **Faltam recursos** e Estratégias Viáveis para essa elaboração.



# Introdução – Justificativa

- Restrição Orçamentária
  - Diminuição de verbas para ciência e tecnologia –2,32%, mesmo com o aumento de base de alunos
  - Aumento do dólar em mais de 3,27% diminuindo o poder de compra
  - Tomada de decisão em saúde, mais de 152mi de Brasileiros dependem exclusivamente do SUS
- Necessidade de dispor estratégias de análise de dados





# Introdução – Justificativa

- Restrição Orçamentária
  - Diminuição de verbas para ciência e tecnologia –2,32%, mesmo com o aumento de base de alunos
  - Aumento do dólar em mais de 3,27% diminuindo o poder de compra
  - Tomada de decisão em saúde, mais de 152mi de Brasileiros dependem exclusivamente do SUS
- Necessidade de dispor estratégias de análise de dados



# Introdução – Abordagem

Utilizar um Cluster Kubernetes como plataforma de orquestração de cargas de trabalho em ambiente virtual.

- Cargas de trabalho:
  - Analise de tendencia de uso de azitromicina entre 2014 e 2021
- Ambientes virtualizados (simulando *Host* do cluster):
  - completa - *Hypervisor* tipo 2
  - sistema operacional - contêineres
- Ambiente virtual:
  - Simulação de máquina de baixo poder computacional:
  - 1 vCPU
  - 2 GB de RAM
  - 6-8 máquinas

Essa Abordagem visa comparar o desempenho desses ambientes simulados, e validar o uso de computadores de baixo poder computacional, no processo de análise de dados de grande volume



# Introdução – Abordagem

O uso de conceitos, métodos e o uso de ferramentas complementares na aplicação da cultura DevOps em ambientes produtivos, permitirá o deployment simplificado melhorando a agilidade e diminuindo a complexidade e operação/sustentação do cluster

- Conceitos como:
  - CI (integração contínua)
  - CD (entrega contínua)
- Uso do método USE (utilização, saturação e erro). Esse método propõe um checklist de métricas a serem coletadas e a avaliação de três parâmetros por meio dessas métricas, relacionando assim o desempenho da carga de trabalho (aplicação) e o desempenho dos nós do cluster sob monitoramento.





# Objetivo

## Objetivos Geral:

Realizar a comparação de desempenho de orquestração de recursos em cluster de baixo custo em ambientes virtualizados, para o processamento e a análise dos dados.

## Objetivos Específicos:

- Realizar a orquestração de recursos em cluster de baixo custo;
- Comparar o desempenho de clusters em ambientes virtualizados;
- Validar o uso de um cluster de utilização compartilhada para processamento de dados distribuídos;
- Propor um método de análise em cluster Kubernetes com uso de computadores desktops;

|





# Revisão de literatura- Análise de dados

- Descisões em saúde costumam ser complexas – precisam de suporte científico (dados) e avaliação de Contexto
- Com o crescimento dos 3V's de dados na área da saúde (Big Data) processar e analisar esses dados tornouse fundamental para tomada de descisões adequadas
- Desafios:
  - complexidade dos dados obtidos
  - ausencia de validação de sistemas, métodos e ferramentas para o tratamento de dados na área
  - custos de novos equipamentos capazes de analisar tal volume
- Há grande oportunidade para a proposição de estratégias de processamento e análise de dados na área



# Revisão de literatura – Alternativas open source

- Considerando
  - O escopo deste trabalho
  - As estratégias para processamento e análise de dados disponíveis no mercado

As soluções encontradas no mercado foram agrupadas em dois grupos:

- Soluções de Computação em nuvem privada:
  - ▶ Se estendem para além do propósito desse trabalho
  - ▶ Requisitos de hardware elevados
  - ▶ Complexidade de configuração devido a sua abrangência

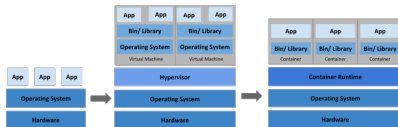
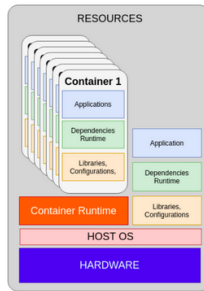




# Revisão de literatura- Alternativas open source

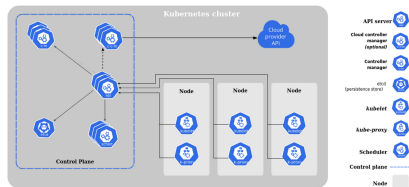
- Soluções de Orquestração de Containers:

- Kubernetes®
- Apache Mesos®
- Hashicorp Nomad®
- Docker Swarm®



# Revisão de literatura- Cluster orquestrador de container

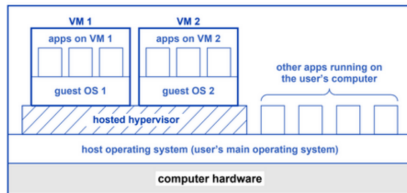
- Kubernetes®:
  - Origem de 15 anos de trabalho da Google (Borg)
  - Estrutura de objetos componentizados
    - ▶ Kube-apiserver
    - ▶ Kube-scheduler
    - ▶ Kube-controller-manager
    - ▶ Kubelet
    - ▶ Kube-proxy
    - ▶ Pod





# Método - Especificação dos nós integrantes cluster de baixo custo

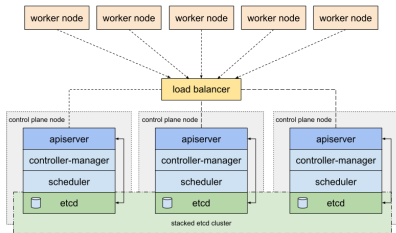
- Cluster Simulado:
  - Virtualização:
    - ▶ Maquinas Virtuais (VMs) (*Hypervisor* tipo 2)
    - ▶ Contêineres Aninhados (Docker In Docker, ou DinD)
  - Especificações de hardware  
1vCPU, 2 GB de RAM;
- provisionamento em 2 etapas
- máquinas subutilizadas
- CAPEX



# Método - Plataforma de orquestração de carga de trabalho

- Arquitetura sugerida para produção:
  - Multi-master com Etcd junto ao nó master
- Alta disponibilidade do cluster
- Recursos de hardware limitados

kubeadm HA topology - stacked etcd



# Método - Plataforma de orquestração de carga de trabalho

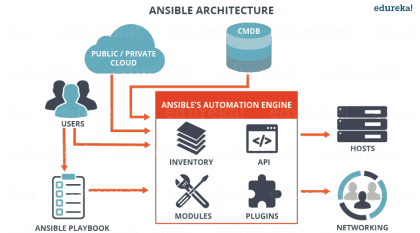
- Implantação da carga de Trabalho
  - Container
  - Parametrizável
  - Volume compartilhado



# Método - Configuração e provisionamento do cluster

O uso de gerenciadores de configuração garantem o versionamento das configurações permitindo maior controle sobre as propriedades dos *assets* gerenciados (Ansible®)

- *Agentless*
- Idempotência
- Gerenciamento de inventário
- SSH - Escolha do algoritmo de criptografia



# Método – Monitoramento

- *OpenTelemetry*
- *Prometheus* Monitoramento de sistemas e Banco de dados de series temporais
- *Grafana* – Dashboard e observabilidade
- Parametros de tempo, taxa de utilização de memoria e processamento

