

Conceitos de Ciência dos Dados e Big Data

PCS5787 – Ciência dos Dados
Prof. Dr. Pedro Luiz Pizzigatti Corrêa
24 de Setembro de 2020

Agenda

- Introdução
 - O conceito de Dado x Informação x Conhecimento
 - Big Data, Ciência dos Dados, Gestão dos Dados
- Gestão dos dados – Importância e Desafios
- O modelo de dados do ponto de vista do *Big Data*
- Desafios relacionados ao *Big Data* e Ciência dos Dados
- Exercícios



Introdução

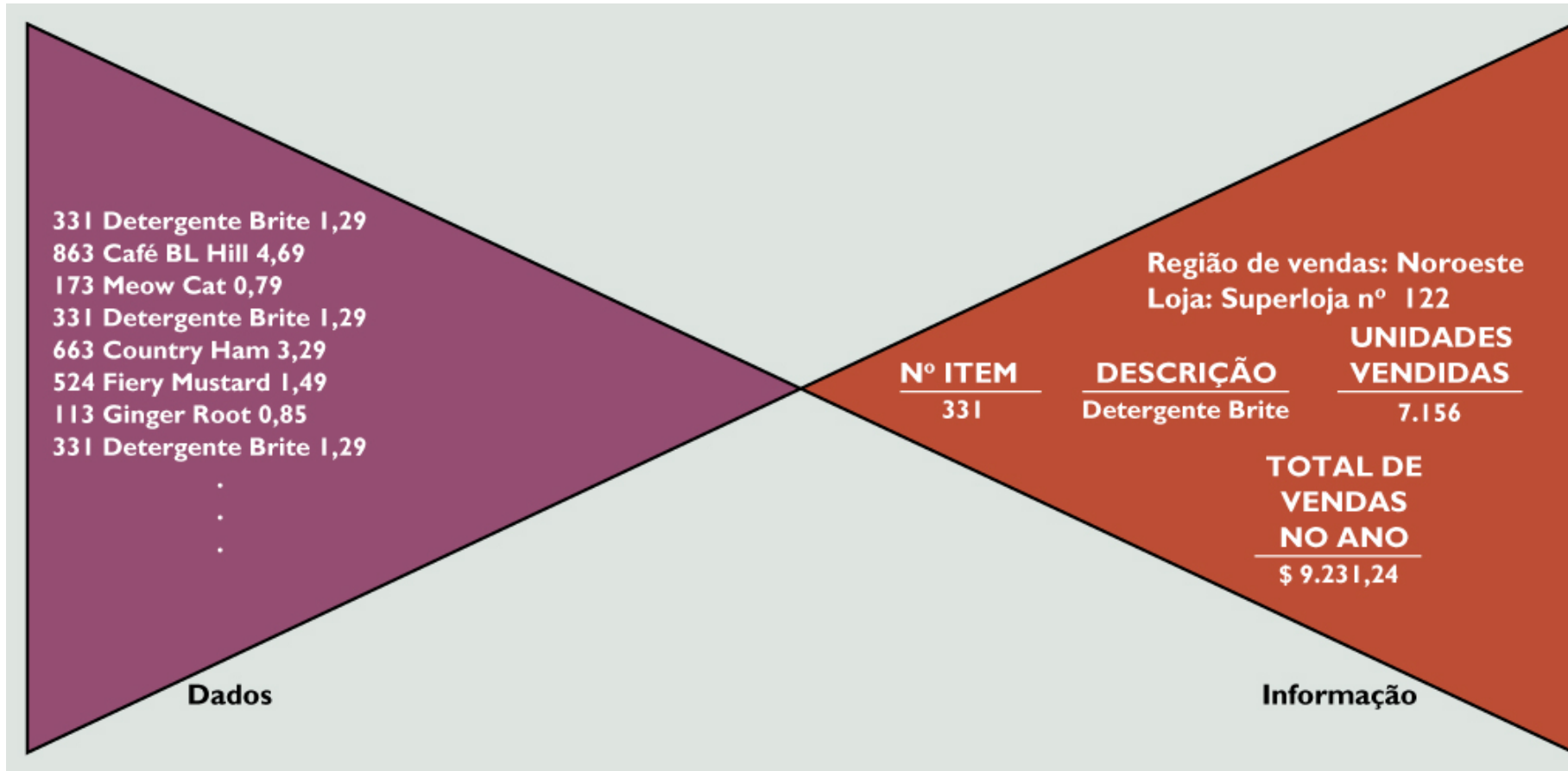
O conceito de Dado x Informação x Conhecimento

CONCEITO DADOS, INFORMAÇÃO E CONHECIMENTO

1. **Dados:** Fluxos de fatos coletados (brutos) que representam eventos do domínio. Qualquer evento que possa ser armazenado em formato digital, incluindo texto, números, imagens, vídeo ou filmes, áudio, software, algoritmos, equações, animações, modelos, simulações, etc.
2. **Informação:** Conjuntos de dados significativos e úteis a seres humanos em processos como o de tomada de decisões;
3. **Conhecimento:** Informações inter-relacionadas não estruturadas de regras que direcionam as tomadas de decisões.

Fonte: CORRÊA, 2011 – Adaptado Laudon, 2013

O conceito de Dado x Informação x Conhecimento



Quais são os dados ?

Coleções de **registros ou medições** que fornecem um registro de evidências do evento observado “... *qualquer informação que possa ser armazenada em formato digital, incluindo texto, números, imagens, vídeo ou filmes, áudio, software, algoritmos, equações, animações, modelos, simulações, etc.* “

Atributos gerais da informação:

- Digital;
- Heterogêneo;
- Contextualizado;
- Valioso.

Cortesia: Profa. Dra. Suzie Allard (University of Tennessee)

O conceito de Dado x Informação x Conhecimento

Pergunta:

Os dados climáticos coletados por estações metrológicas numa região é dado ou informação ?



BIG WORLD

BIG PROBLEMS

BIG DATA



The
**FOURTH
PARADIGM**

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

Paradigmas da ciência

Há mil anos:

- A ciência foi **empírica**.
- Usada para descrever fenômenos naturais.



Observações

Há poucos séculos:

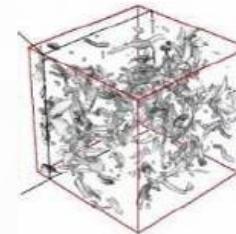
- A ciência passou a ser também **teórica**.
- Uso de modelos, generalizações, etc.

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$

Leis de Kepler, Newton, Maxwell

Nas últimas décadas:

- Pesquisadores passaram a validar seus modelos teóricos com o uso de simulações.
- Ciência **computacional**.



Simulação de fenômenos complexos

e-Science: O quarto paradigma

Hoje:

Ciência orientada a grande volume de dados

(*Data-intensive Science*: Unifica teoria, experimentação e simulação).

- Dados capturados por instrumentos ou gerados por simulação.
- Dados processados por software.
- Informação/conhecimento armazenados em computadores (**em grande escala**).
- Pesquisadores analisam arquivos/bases de dados por meio de gerenciamento de dados e estatísticas.

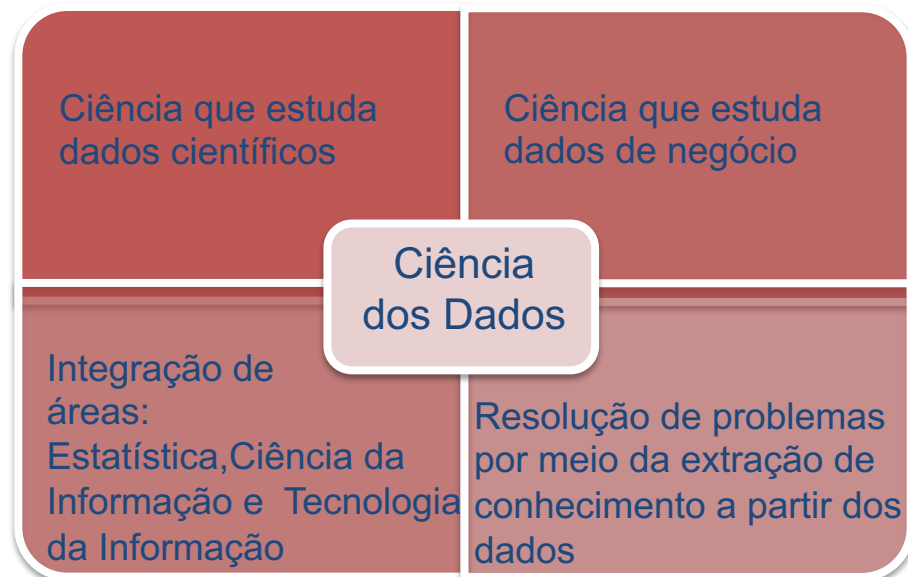
Três atividades consideradas na exploração de dados:

- Captura
- Curadoria
- Análise



Ciência dos Dados

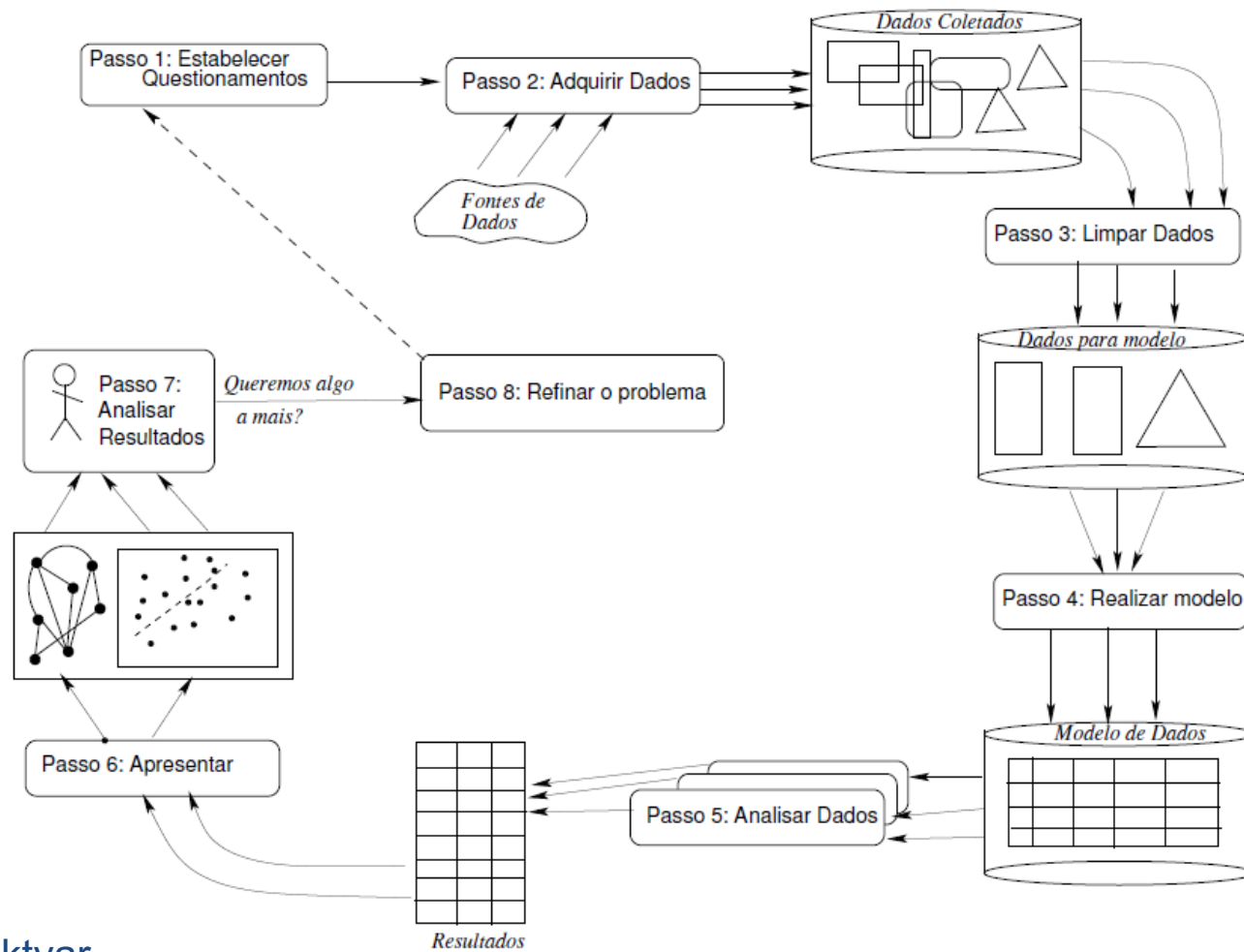
- Em busca por uma definição formal sobre Ciência dos Dados, encontramos diversos trabalhos na literatura
 - Embora muito se discuta sobre a composição das atividades de Ciência dos Dados, o seu conceito ainda não é algo fundamentalmente estabelecido
- Para Zhu e Xiong (2015), há quatro vertentes (perspectivas) que buscam caracterizar Ciência dos Dados



Ciência dos Dados

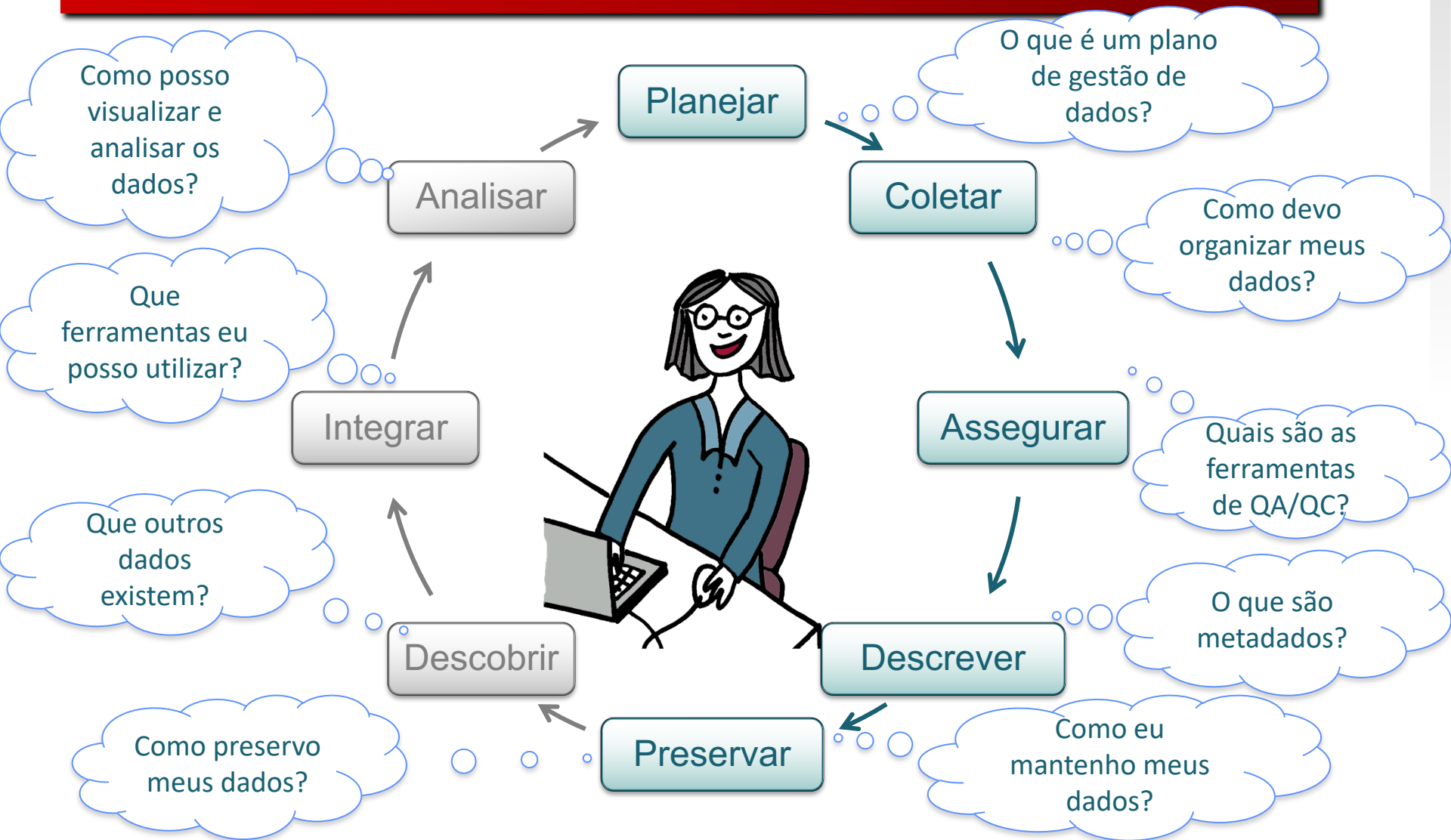
- Embora não haja consenso sobre a definição, encontramos como elemento comum em todas as propostas um processo de manipulação, processamento e análise de dados, que visa a descoberta de novos conhecimentos
- Para Alex Dehktyar (2016),
 - Ciência dos dados é uma disciplina que permite tratar o ciclo de trabalho com os dados, considerando atividades que compreendem desde a aquisição dos dados, passando pela análise dos dados, até o processo de apresentação dos dados e obtenção de novos conhecimentos

Ciência dos Dados - Processo



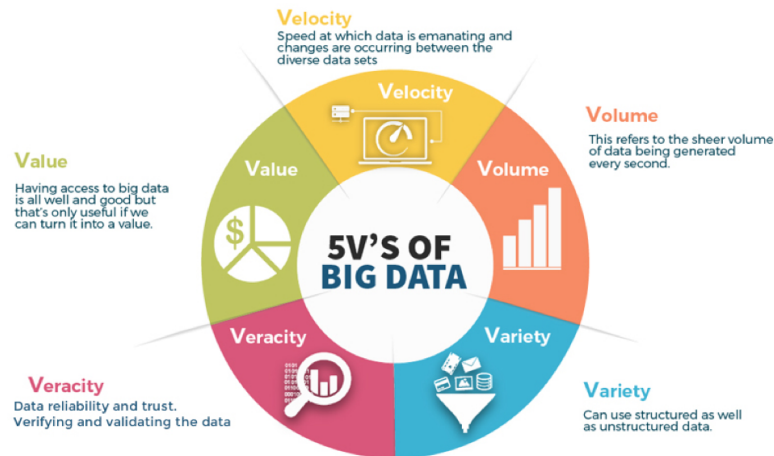
Cortesia: Alex Dehkyar

Ciência dos Dados - Gestão de Dados

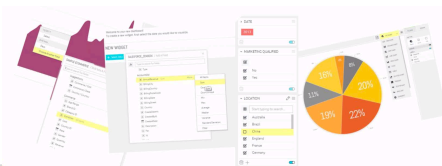


O que é Big Data ?

- ✓ É uma aplicação computacional de Ciência dos Dados que tem por objetivo analisar, extrair sistematicamente informações de grandes volumes de conjuntos de dados, para os quais técnicas computacionais tradicionais não são adequadas. Os desafios para gestão dos dados são classificados em 5V's (Chen et al., 2012, Kwon et al., 2014).



- ✓ Big data é um grande volume de dados, alta velocidade e alta variedade de ativos de informação que demandam formas inovadoras e econômicas de processamento de informações para melhor insight e tomada de decisões.” (“Gartner IT Glossary, n.d.”)



Causas que tornam os dados complexos

- Linguagem de consulta
- Tamanho
- Estrutura
- Dispersão
- Taxa de crescimento



Etapas para criação de Big Data

Questões a tratadas:

Dispersão

Estrutura

**Tamanho e taxa
de crescimento**

**Linguagem de
consulta e
detalhes**



Extrai
Transforma
carrega
atualiza



Consultas

Análise
Consulta/
Relatório

Etapas:

Fontes de Dados ETL

Centralização

Análise

Apoio a tomada de decisões

- Apoio para tomar decisão:
 - *Quais foram os volumes de vendas por região e categoria de produto no último ano ?*
 - *Quais pedidos devem ser priorizados para maximizar os lucros ?*
 - *10% de desconto aumentará o volume de vendas suficientemente?*
- On-Line Analytical Processing (OLAP) é uma ferramenta de Sistemas de Suporte a Decisão (Decision Support Systems - DSS) e Sistema de Apoio a Executivos (SAE ou Executive Information System -EIS)

Evolução:

- 60: Relatórios Batches
 - Longos relatórios;
 - Inflexíveis e caros, necessitando reprogramar a cada nova requisição
- 70: DSS e EIS usando terminais
 - Ainda Inflexível;
 - Não integrados a outras ferramentas de *desktop*.
- 80: Acesso através de *Desktop* e ferramentas de análise:
 - Ferramentas de consultas (*queries*), planilhas e GUI;
 - Fácil de usar, mas acessa somente a bases de dados operacionais
- 90: *Data warehousing* com OLAP e outras ferramentas integradas.
- 00/10: Aplicações/Serviços Web e Big Data

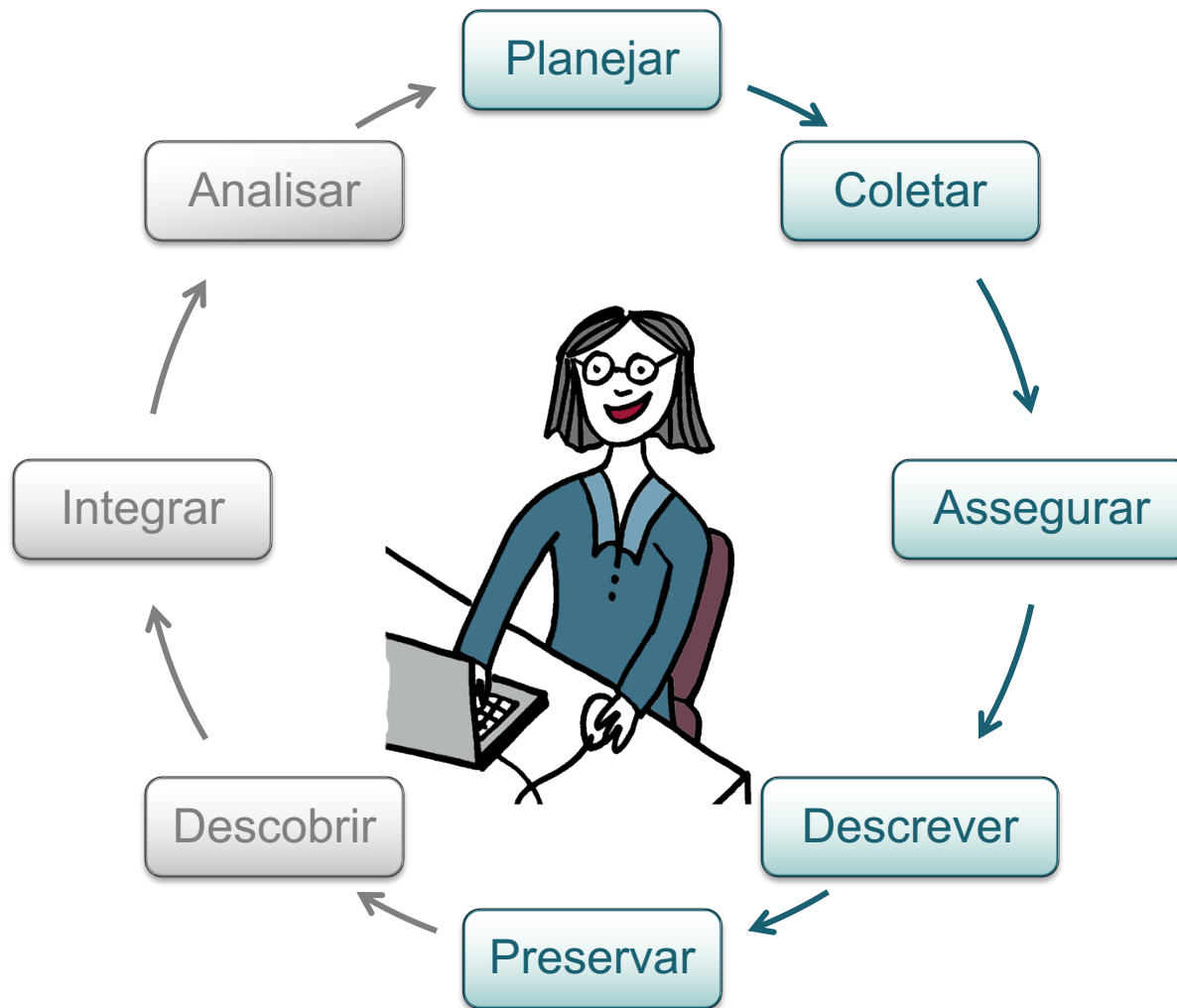
Agenda

- Introdução
 - O conceito de Dado x Informação x Conhecimento
 - Big Data, Ciência dos Dados, Gestão dos Dados
- Gestão dos dados – Importância e Desafios
- O modelo de dados do ponto de vista do *Big Data*
- Desafios relacionados ao *Big Data*
- Exercícios



Gestão de Dados – Importância e Desafios

Gestão de Dados



2. Por que a gestão de dados?

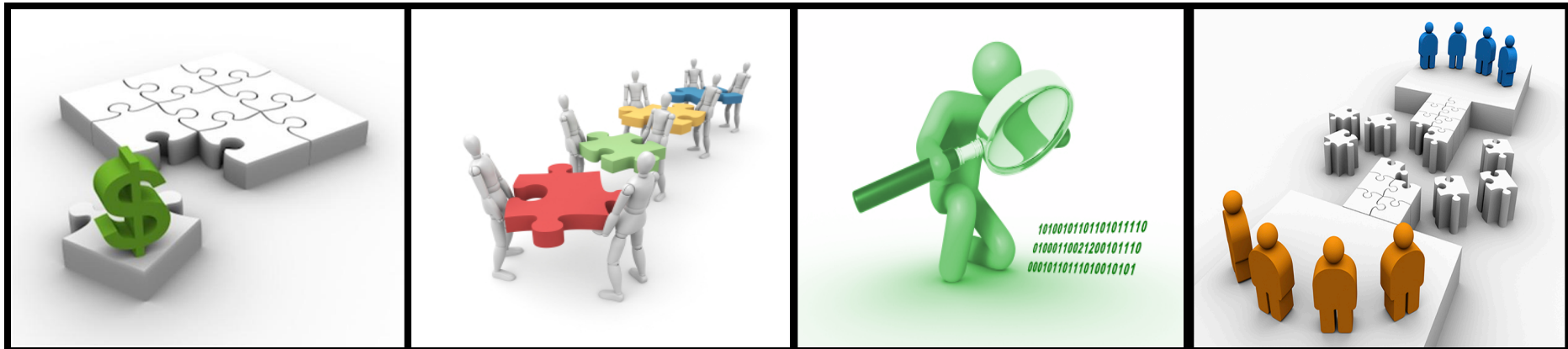


PORQUE APLICAR AS TÉCNICAS E CONCEITOS DE GESTÃO DE DADOS?

2. Porque a gestão de dados?

1. Para capturar, armazenar, proteger e garantir a integridade dos ativos de dados;
2. Garantir a utilização adequada dos dados e informações;
3. Maximizar o uso eficaz dos dados e agregar valor aos ativos da informação.

Fonte: DAMA International, *The DAMA Guide to the Data Management Body of Knowledge*



2. Porque a gestão de dados?

Se seus dados caírem em mãos erradas?

national security has leaked from Whitehall, the head of the civil service has warned.

Cabinet Secretary Sir Gus O'Donnell said there were "one or two" leaks from areas dealing with national security.

However he said that "leaking sensitive official information have hit the headlines

There has been a series of breaches in recent years, including of the entire child benefit records, with the personal details of 25 million people.

Giving evidence to the Commons Public Administration Committee, Sir



INFORMAÇÕES COM POTENCIAIS IMPLICAÇÕES PARA A SEGURANÇA NACIONAL VAZARAM DO GOVERNO BRITÂNICO.

Fonte:

<http://news.bbc.co.uk/1/hi/uk/8332445.stm>

2. Por que a gestão de dados?

SE FOR NECESSÁRIA REPRODUZIR AS ANÁLISES?

change row

Leading British scientists at the University of East Anglia, who were accused of manipulating climate change data - dubbed Climategate - have agreed to publish their figures in full.



CIENTISTAS FORAM ACUSADOS DE MANIPULAR DADOS SOBRE MUDANÇAS CLIMÁTICAS.

Print this article

Twitter 3

Email

LinkedIn 0

Copenhagen climate change conference
News » UK News »
Earth News »

Fonte: The
Telegraph

2. Por que a gestão de dados?

SE ESTE FOR O SEU INSTITUTO DE PESQUISA?



Incêndio no Instituto Butantan destrói maior acervo de cobras do país

Fogo queimou 70 mil espécies conservadas em formol na Zona Oeste de SP. Chamas atingiram laboratório de répteis; causas do fogo serão apuradas.


Fonte: <http://g1.globo.com/sao-paulo/noticia/2010/05/incendio-no-instituto-butantan-destroi-maior-acervo-de-cobras-do-pais.html>

2. Por que a gestão de dados?

SE ESTA FOR A SUA
MOCHILA?

“O HD externo é muito importante, pois contém 5 anos de backup de e-mails”

CASH REWARD
for returning my lost backpack



- Black [AK] Burton Rucksack
- Lost on Friday 15. July at 8 pm in the Panton Arms pub 43, Panton St. Cambridge
- Containing a laptop (white MacBook), a black external hard drive and scientific research documents

The external hard drive is VERY important to me as it contains 5 years of research data which are crucial for my PhD thesis!!!
If you found it, I would be extremely grateful if you could return it to the Panton Arms or contact me on: 07804430054 (ar456@cam.ac.uk)

Thank you!!

Fonte: <http://blogs.ch.cam.ac.uk/pmr/2011/08/01/why-you-need-a-data-management-plan>

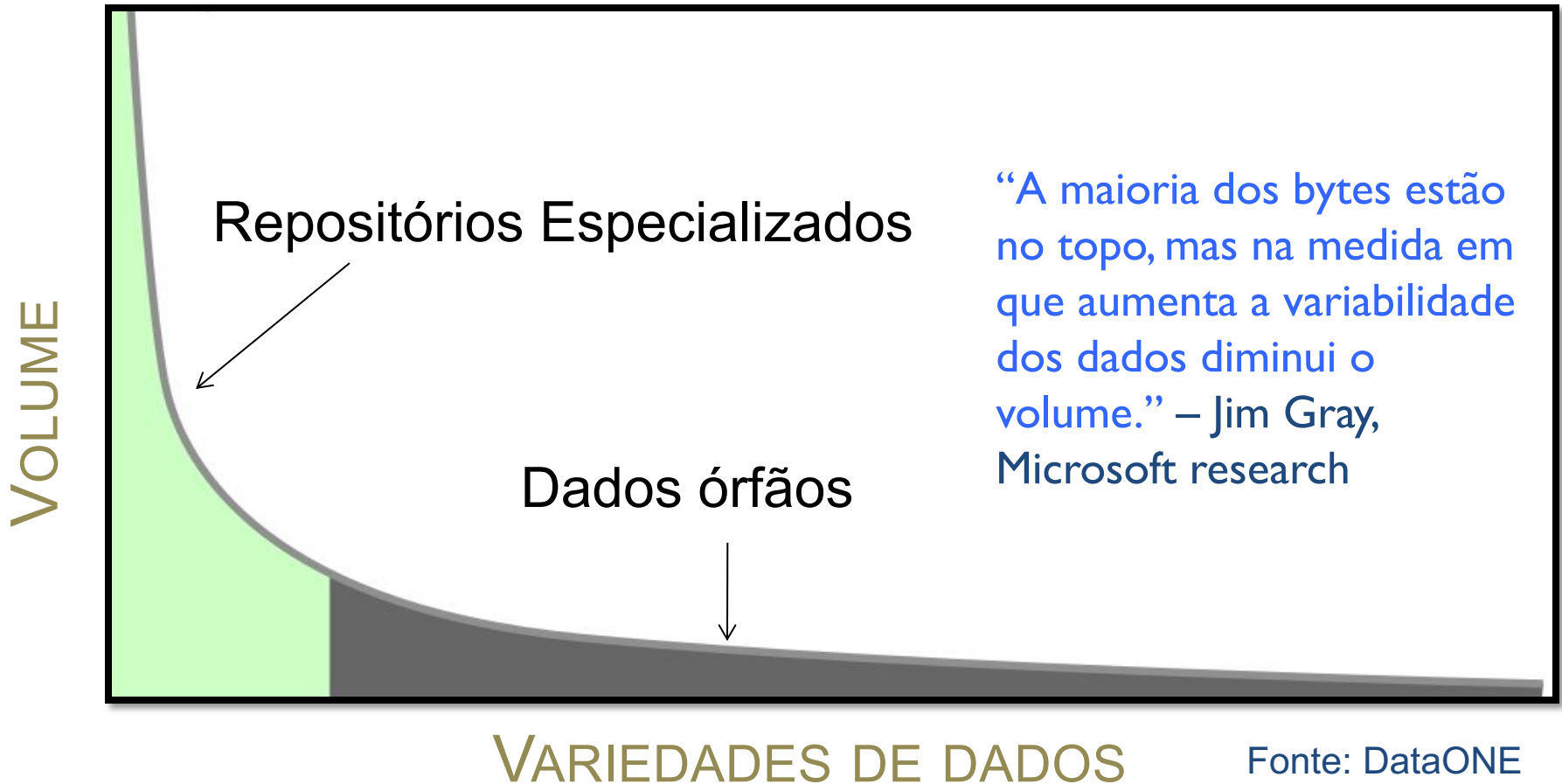
2. Gestão de dados

“Gestão de Dados é a disciplina responsável por definir, planejar, implantar e executar: estratégias, procedimentos e práticas necessárias para gerenciar de forma efetiva os recursos de dados e informações das organizações, incluindo planos para sua definição, padronização, organização, proteção e utilização.”

Fonte: DAMA-DMBOK

A Gestão de Dados é um conceito bastante amplo, ela atua nos níveis: Operacional, Gerencial (Tática) e Estratégico.

2. Desafios: “The Long tail” da Gestão dos Dados



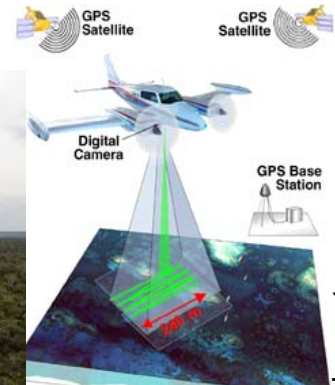
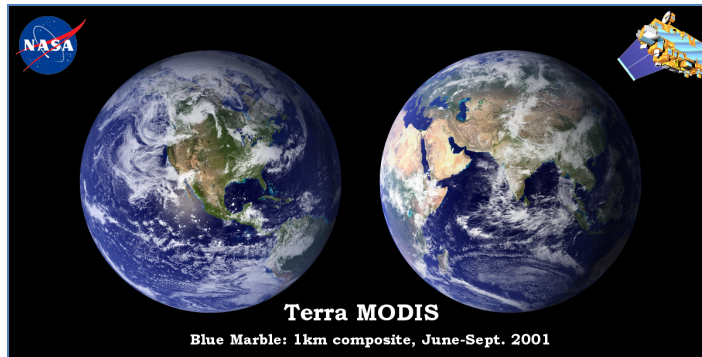
2. Desafios: Dados órfãos

- INFORMAÇÃO QUE SE TORNOU IRRECUPERÁVEL POR ESTAR LOCALIZADA EM DISPOSITIVOS NÃO MAIS ACESSÍVEIS, COMO NOTEBOOKS, E QUE NUNCA FORAM TRANSFERIDAS PARA SERVIDORES COMPUTACIONAIS;
- INFORMAÇÕES PERDIDAS APÓS O DESLIGAMENTO DE PESQUISADORES/FUNCIÓNÁRIOS DA INSTITUIÇÃO;
- DADOS DE PESQUISADORES NÃO ASSOCIADOS A NENHUMA REDE DE DADOS.



2.Desafios: “Dilúvio” dos Dados

Redes, Sensores, Sensoriamento Remoto, Experimentos, Coletas...

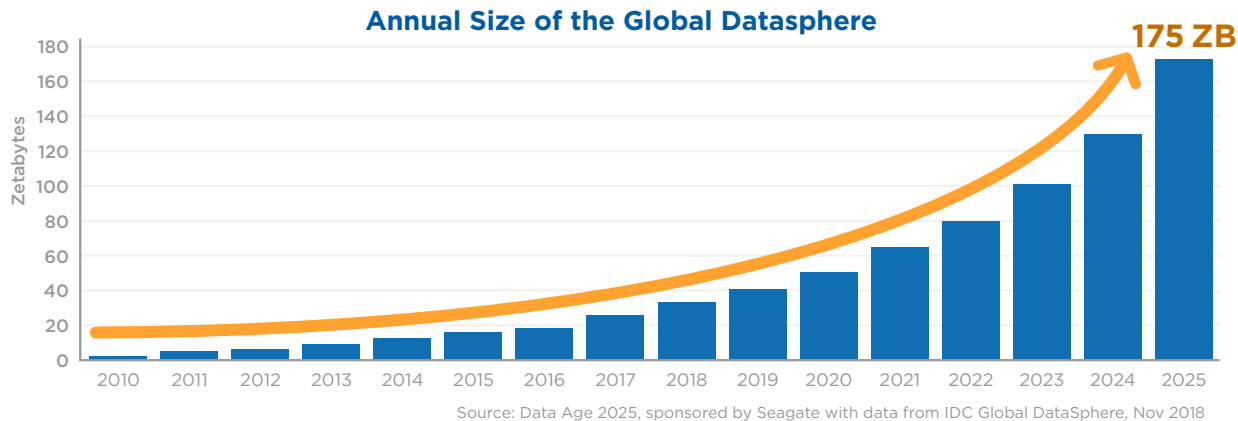


Fonte: www.carboafrika.net

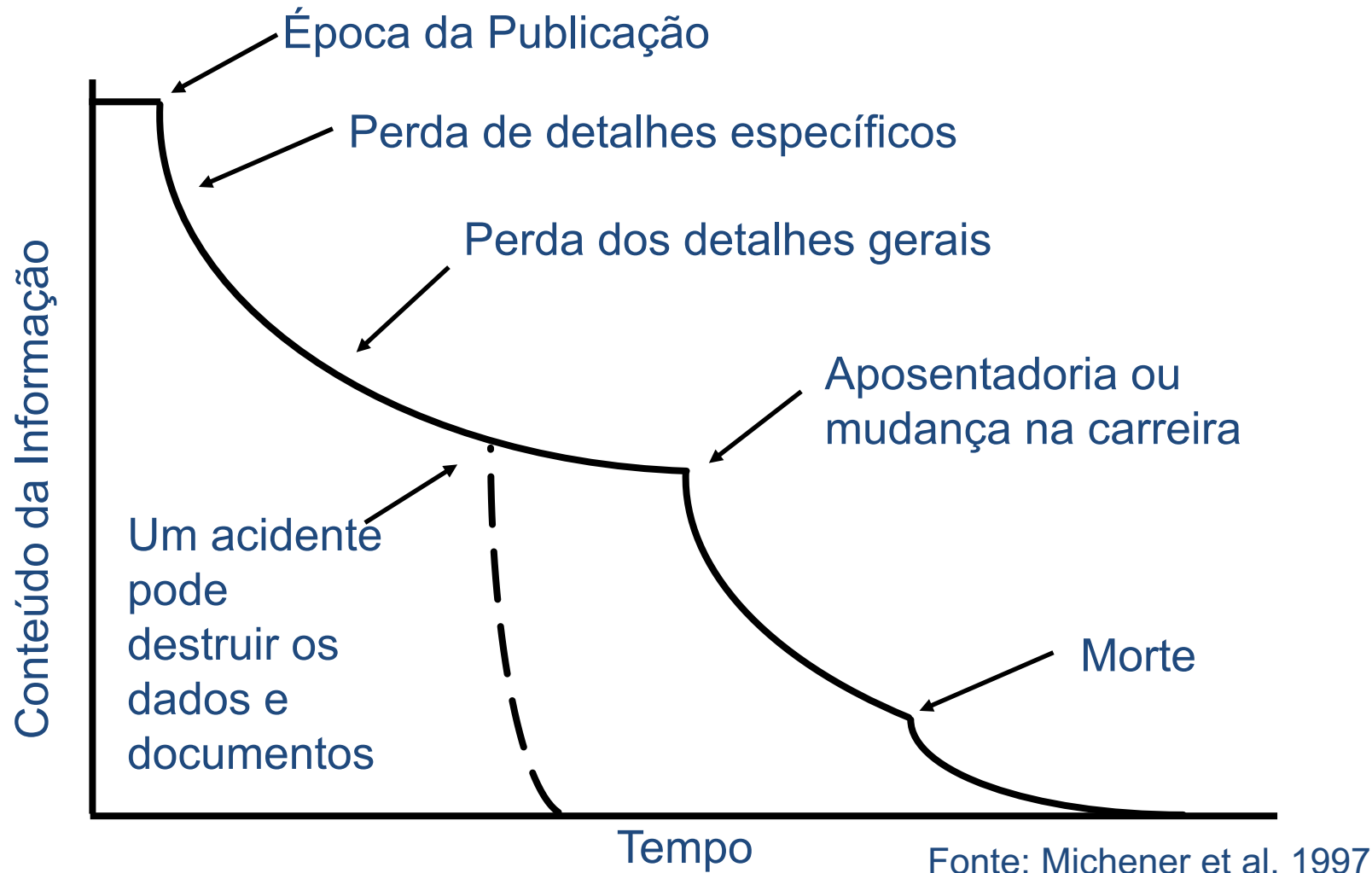
Problema: Volume de Dados

International Data Corporation (IDC):

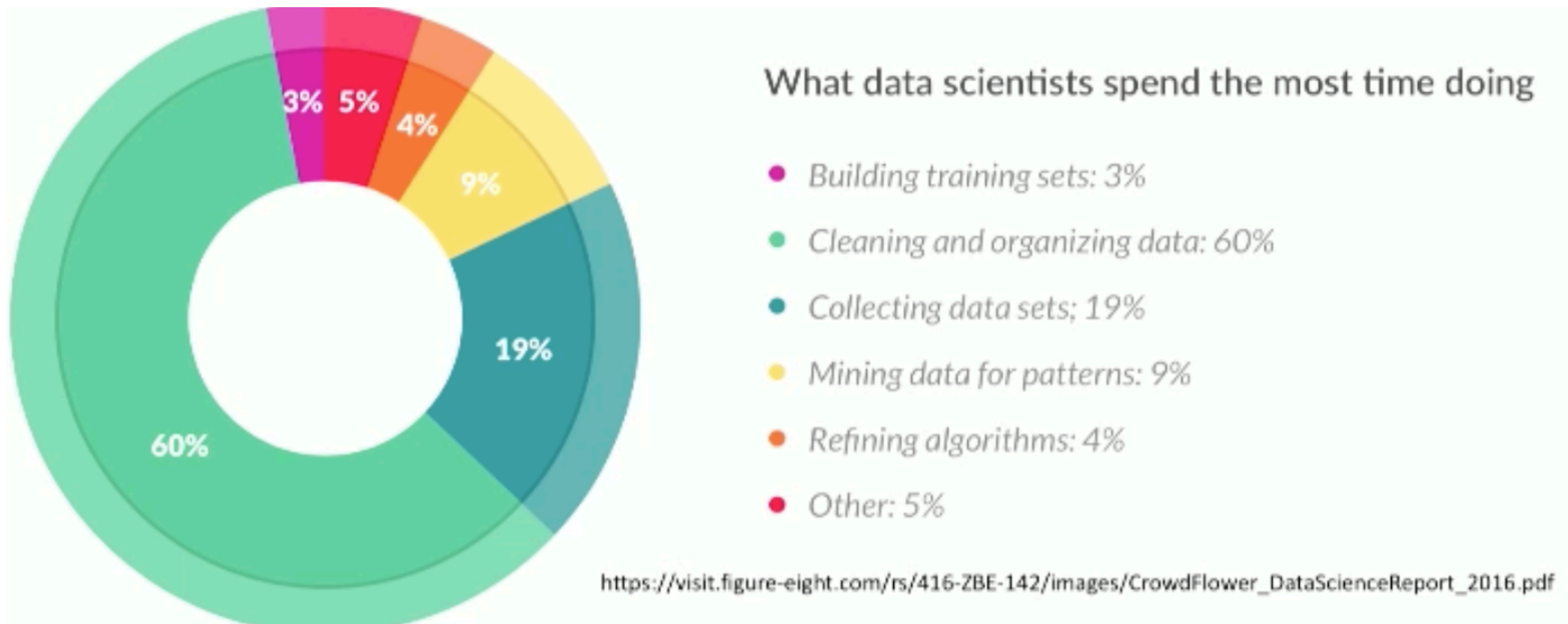
- A informação no mundo está mais que dobrando, está dobrando cada dois anos (30 Zettabytes em 2018);
- O número de arquivos está crescendo mais rápido que a capacidade de armazenamento. Nos próximos 5 anos esses arquivos irão crescer num fator de 8x;
- O número de pessoas na área de IT responsável por esses dados irá crescer “suavemente”;



2.Desafios: Entropia dos Dados



Problema que conhecemos muito bem !



- 60% em atividades na organização e limpeza dos dados;
- 19% em atividades de pré-processamento
- 9% em atividades de análise

Agenda

- Introdução
 - O conceito de Dado x Informação x Conhecimento
 - Big Data, Ciência dos Dados, Gestão dos Dados
- Gestão dos dados – Importância e Desafios
- O modelo de dados do ponto de vista do *Big Data*
- Desafios relacionados ao *Big Data*
- Exercícios



Gestão de Dados – Boas Práticas e Modelo

2. BOAS PRÁTICAS NA GESTÃO DE DADOS

Planejar



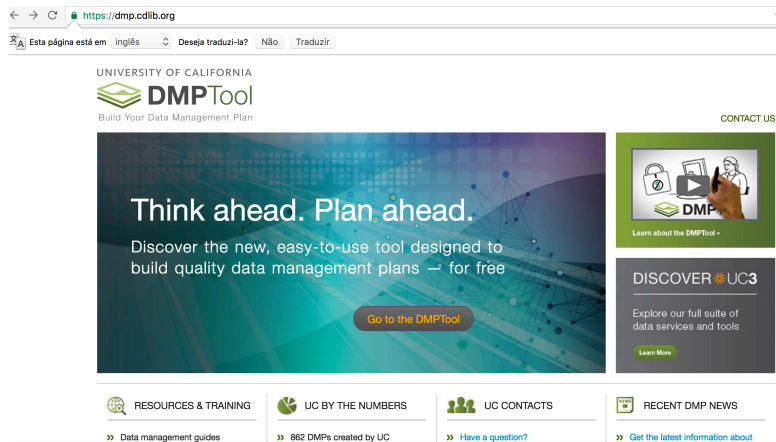
Coletar

2 - Planejar

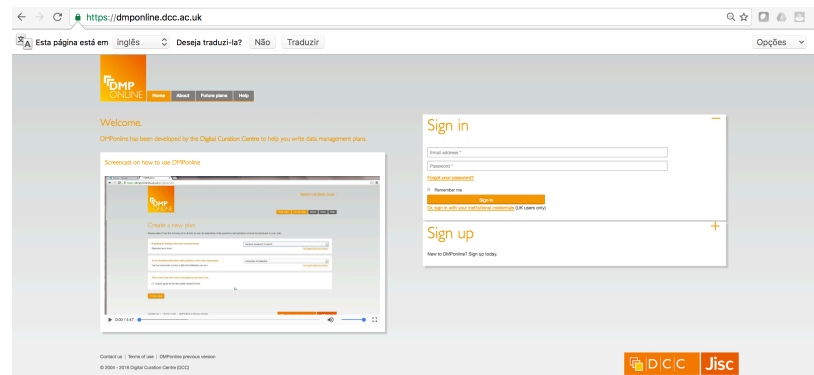
- a) Mapeia os processos e recursos para todo o ciclo de vida dos dados;
- Começar com os objetivos do projeto (desejos, resultados e impactos) e construir um plano de gestão dos dados, considerando a política dos dados e sua sustentabilidade.
 - Pontos a serem considerados:
 - Quais e como os dados serão Coletados;
 - Onde os dados serão armazenados (Repositório);
 - Como os dados serão organizados (formatos, estrutura);
 - Como os dados serão descritos (metadados);
 - Como os dados serão compartilhados;
 - Definição do plano de preservação dos dados e responsabilidades.

2 - Plano de Gestão dos Dados

b) Ferramentas para a criação de Planos de Gestão de Dados:




dmp.cdlib.org



dmponline.dcc.ac.uk

2 - Plano de Gestão dos Dados

← → ↻ <https://dmptool.org> ☆

 **DMPTool**

[Home](#) [DMP Requirements](#) [Public DMPs](#) [News](#) [Help](#) [Contact Us](#) [About](#) [Log In](#)

Data Management Planning Tool

Create, review, and share data management plans that meet institutional and funder requirements.

[Get Started](#)



PUBLIC DMPs

List of sample data management plans provided by DMPTool users.



DMPTOOL NEWS

Latest information about data management and the DMPTool.



DMPTOOL HELP

Overview of how to use the tool, plus resources and guidance on data management.

BOAS PRÁTICAS NA GESTÃO DE DADOS

Planejar



Coletar



Assegurar

2 - Coletar

1. **Determina a melhor forma para se obter os dados;**
2. **O resultado deste processo é um Documento que descreve a forma como os dados são estruturados.**
 - **Determina como os dados são coletados;**
 - **Determina como os dados gerados são organizados e armazenados.**
3. **Para o armazenamento em longo prazo, os dados devem ser gerados em formatos de dados consistentes e específicos para esta finalidade (facilitando sua utilização hoje e no futuro).**
 - **Bases Relacionais;**
 - **Arquivos CSV (Texto com separadores);**
 - **Outros recomendações de formatos**

<https://libraries.mit.edu/data-management/store/formats/>

Problemas na coleta

1	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Site	Date	Plot	Species	Weight	Acult		Rodent Trapping 3/15/2010						
2	DeepWell	2/13/2010	1	DIPO	12.1	j		Site	Plot	Adult	RodentSp	Weight		
3	Deep Well	Feb-10	2	Pero	13.22	j		DW		1	y	Pero	12	
4	rioSalado	2/13/2010	1a	pero	16	N		RS		2	j	PERO	escaped <15	
5	riuSladu	"	1*	CleGap	18.92	gut away		RS		3	ri	Clegap	91	
6				Mean1	15.06									
7														
8														
9														
10														
11														
12	Rodent Trapping			MJK & ALN	10-Apr-10									
13	Site	Plot	Adult	Species	grams	Cmments								
14	deep well		1	y	woodrat	13								
15	riosalado		2	y	PERO	24.5								
16	riosalado		3	y	Clegap	91								
17														
18														
19														
20														

- Inconsistências entre os eventos coletados;
- Localização das informações de datas;
- Inconsistência no formato de datas;
- Nome das colunas;
- Ordem das colunas.

Problemas na coleta

A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	Site	Date	Plot	Species	Weight	Adult			Rodent Trapping	3/15/2010			
2	DeepWell	2/13/2010		1 DIPO	12.1	j		Site	Plot	Adult	RodentSp	Weight	
3	Deep Well	Feb-10		2 Pero	13.22	j		DW		1 y	Pero		12
4	rioSalado	2/13/2010	1a	pero	16	N		RS		2 j	PERO	escaped <15	
5	riuSladu	"	1*	CleGap	18.92		gut away	RS		3 n	Clegap		91
6				Mean1	15.06								
7													
8													
9													
10													
11													
12	Rodent Trapping		MJK & ALN	10-Apr-10									
13	Site	Plot	Adult	Species	grams		Comments						
14	deep well		1 y	woodrat		13							
15	riosalado		2 y	PERO		24.5							
16	riosalado		3 y	Clegap		91							
17													
18													
19													
20													

	A	B	C	D	E	F	G	H
1	Date	Site	Plot	Species	Weight	Adult	Comments	
2	2/5/2010	Deep Well		1 DIPO	13.2	y		
3	2/4/2010	Deep Well		1 CLEGAP	11.6	j		
4	2/5/2010	Rio Salado		1 DIPO	14.2	y		
5	2/5/2010	Rio Salado		2 PERO	10.1	y		
6	3/15/2010	Deep Well		1 DIPO	15.2	y	plot burned	
7	3/15/2010	Deep Well		2 DIPO	21.7	y	pregnant	
8	3/15/2010	Rio Salado		1 CLEGAP	16.2	j		
9								
10								
11								
12								
13								

- Colunas devem ser consistentes: apenas números, datas ou textos;
- Formatos e códigos consistentes;
- Dados em apenas uma planilha, para facilitar a captura dos dados por sistemas computacionais, sem a intervenção humana;
- Colunas específicas para comentários e demais características que descrevem o registro.

BOAS PRÁTICAS NA GESTÃO DE DADOS

Coletar



Assegurar



Descrever

3 - Assegurar

1. Empregar procedimentos de garantia e controle da qualidade (QA/QC) dos dados:
 - Capacitação de Coletores;
 - Rotina de calibração de instrumentos;
 - Procedimentos de revisão dos dados.
2. Identificar problemas e técnicas computacionais possíveis para solucioná-los.

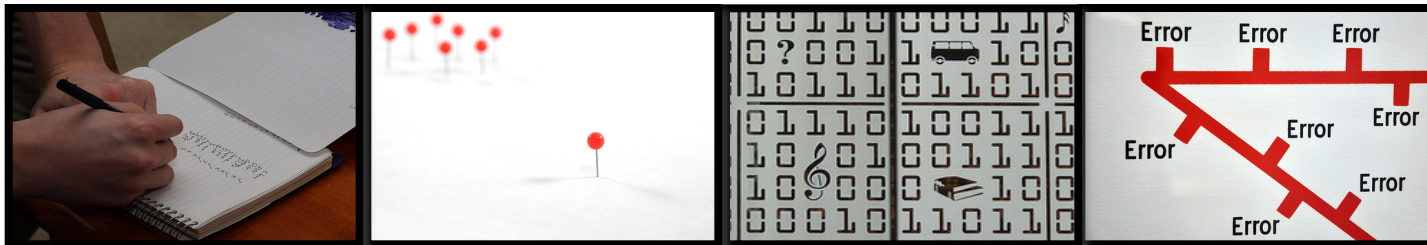
Quality Assurance (QA) é um conjunto de processos utilizados para garantir que os dados necessários serão coletados e armazenados;

Quality Control (QC) é um conjunto de processos para avaliar a qualidade dos dados depois que os mesmos forem coletados.

3-Assegurar

1. Tipos mais comuns de erros:

- Erros de Ação
 - Dados inseridos incorretamente;
 - Exemplos: Instrumentos mal calibrados, erros de digitação.
- Erros de Omissão
 - Dados ou Metadados não registrados;
 - Exemplos: Documentação inadequada, erro humano, anomalias nas coletas.



BOAS PRÁTICAS NA GESTÃO DE DADOS

Assegurar



Descrever

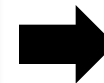
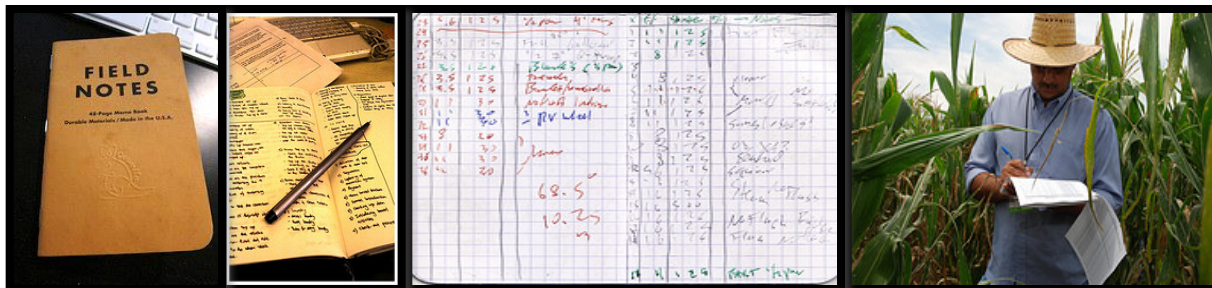


Preservar

4 - Descrever

1. Documentação dos dados

- Porque, quem, o quê, quando, onde e como;
- Utilização de Metadados (dados sobre dados);
 - É a **chave** para o compartilhamento e reutilização dos dados.
 - Várias normas e ferramentas estão disponíveis para apoiar este processo.
 - Os metadados são utilizados tanto por humanos quanto por computadores para apoio nos processos de descoberta, integração e análise dos dados.



4- Descrever

1. Reutilização dos dados

- Se são criados metadados, outras pessoas podem descobrir seus dados.
- Se são criados metadados, você pode encontrar os seus próprios dados!



BOAS PRÁTICAS NA GESTÃO DE DADOS

Descrever



Preservar

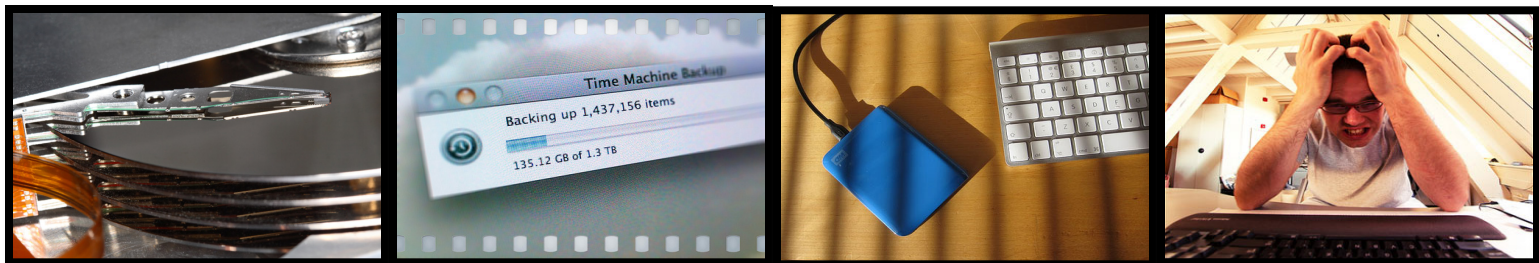


Descobrir

5 - Preservar

Plano para preservar os dados

- No curto prazo para minimizar as perdas potenciais, como acidentes;
- No Longo prazo para que participantes do projeto e outros pesquisadores possam acessar, interpretar e utilizar os dados no futuro.
- Proteção dos dados incluem questões como backups, segurança física, criptografia;
- Preservação dos dados incluem os processos para conservação, recuperação, reorganização e descrição dos dados.



5 – Preservar

Boas práticas

- **Armazenar os dados e metadados em formatos apropriados;**
- **Realizar backups dos dados dos projetos;**
- **Garantir a integridade e o acesso aos backups de dados;**
- **Definição de políticas para segurança e armazenamento das informações;**
- **Identificar dados com valor a longo prazo;**
- **Preservar os dados brutos;**
- **Identificar e gerenciar os dados sensíveis.**

5 – Preservar - exemplo

- DDOMP Checklist – Team Resources

- Material development and temporary storage location
 - Google Drive
- Team communications and information decimation tools
 - Email, Slack
- Dataset storage location during the project
 - Open Science Framework (AWS integration)
- Software development platform
 - GitHub
- Data preservation (including derived products) repository
 - Environmental Data Initiative
- Software preservation repository
 - Zenodo
- Training, workshop material preservation repository
 - Zenodo

- Stall, S. DDOMP - [10.5281/zenodo.3942688](https://zenodo.org/record/3942688)
- Projeto Parsec: <https://parsecproject.org/>

PARSEC

- PIs – 4
- Country Leaders – 6
- Funders – 4
- Researchers – 30
- Languages - 4

5 – Preservar - exemplo

- DDOMP Checklist – Tracking and Reporting
 - **Once (during lifetime of researcher)**
 - **ORCID profile** - activate the automatic updates from Crossref (published papers) and DataCite (published datasets and other digital objects. Reference this page for instructions: http://bit.ly/ORCID_Trust)
 - **Weekly**
 - Track **datasets created**, track **datasets used**, track **workflow/provenance**
 - **Monthly**
 - Publish and Report **conference presentations and posters**
 - Peer-reviewed Papers – and Supporting Digital Objects
 - Deposit and preserve **Datasets**
 - Deposit and preserve **Software**
 - Report **publications with citations to datasets, software, and other digital objects**
 - **Quarterly**
 - Update your ORCID profile and ensure accurate and complete to ensure proper credit
- Stall, S. DDOMP - [10.5281/zenodo.3942688](https://zenodo.org/record/3942688)
- Projeto Parsec: <https://parsecproject.org/>

5 – Preservar - exemplo

- PARSEC Data and Digital Output Management Plan and Workbook

Further details can be found on the process and methods used for PARSEC:

Stall, Shelley, Specht, Alison, Corrêa, Pedro Luiz Pizzigatti, David, Romain, Edmunds, Rorie, Mabile, Laurence, Machicao, Jeaneth, O'Brien, Margaret, Wyborn, Lesley. (2020). PARSEC Data and Digital Output Management Plan and Workbook. Zenodo.

[10.5281/zenodo.3891426](https://doi.org/10.5281/zenodo.3891426)

Use your DMP or DDOMP to make your own Checklist.

BOAS PRÁTICAS NA GESTÃO DE DADOS

Preservar



Descobrir



Integrar

6 - Descobrir

1. Estratégias para a localização e aquisição de dados potencialmente úteis:

- Identificar dados complementares que possam agregar valor aos dados do projeto.
- Buscadores genéricos na Web não são eficientes para encontrar dados úteis.
- Dados são mais facilmente encontrados por meio das redes, diretórios de projetos e repositórios:

Repositórios	Website
Global Biodiversity Information Facility	gbif.org
Atlas of Living Australia	ala.org.au
Knowledge Network for Biocomplexity	knb.ecoinformatics.org
Dryad	datadryad.org
DataONE	search.dataone.org

O COMPARTILHAMENTO AGREGA VALOR AOS DADOS

OS METADADOS GARANTEM A PROPRIEDADE DOS DADOS, SUA CONFIABILIDADE E USABILIDADE

OS PATROCINADORES ESPERAM, E ALGUNS EXIGEM, QUE OS DADOS SEJAM COMPARTILHADOS

O COMPARTILHAMENTO DE CONHECIMENTO É ESSENCIAL PARA O AVANÇO DA CIÊNCIA

BOAS PRÁTICAS NA GESTÃO DE DADOS

Descobrir



Integrar



Analisar

FAIR Data Principles

to be Findable, Accessible, Interoperable and Re-Usable



Article

A Data Quality Strategy to Enable FAIR, Programmatic Access across Large, Diverse Data Collections for High Performance Data Analysis

Ben Evans , Kelsey Davies, Jingbo Wang ^{*}, Rai Yang, Clare Richards and Lesley Wyborn

National Computational Infrastructure, for Australian National University, Acton 2011, Australia; Ben.Evans@nicta.edu.au (B.E.); Kelsey.Davies@nicta.edu.au (K.D.); Rai.Yang@nicta.edu.au (R.Y.); Clare.Richards@nicta.edu.au (C.R.); Lesley.Wyborn@nicta.edu.au (L.W.)

^{*} Correspondence: Jingbo.Wang@nicta.edu.au; Tel.: +61-02-4125-0862

Academic Editors: Moushi Gu and Vladimir Dobral

Received: 31 August 2017; Accepted: 8 December 2017; Published: 13 December 2017

Abstract: To ensure seamless, programmatic access to data for High Performance Computing and analysis across multiple research domains, it is vital to have a methodology for standardisation of both data and services. At the Australian National Computational Infrastructure (NCI) we developed a Data Quality Strategy (DQS) that currently provides processes for: (1) Control of data structures needed for a High Performance Data (HPD) platform; (2) Quality Control through compliance with recognised community standards; (3) Benchmarking cases of open performance tests; and (4) Quality Assurance (QA) of data through demonstrated format and performance across common platforms, tools and services. By implementing the NCI we have seen progressive improvement in the quality and usefulness of the datasets across different subject domains, and demonstrated the ease by which modern programmatic methods can be used to access the data, either in situ or via web services, and for uses ranging from traditional analysis methods through to emerging machine learning techniques. To help increase data re-use by broader communities, particularly in high performance environments, the DQS is also identifying the need for any extensions to the relevant international standards for interoperability and programmatic access.

Keywords: data quality; quality control; quality assurance; benchmarks; performance; data management policy; netCDF; high performance computing; HPC; big data

1. Introduction

The National Computational Infrastructure (NCI) manages one of Australia's largest and diverse repositories (10+ PB) of research data collections spanning datasets from climate, oceans and geophysics through to astronomy, bioinformatics and the social sciences [1]. Within domains, data can be of different types such as gridded, ungridded (i.e., line surveys, point cloud and raster image types), as well as having diverse coordinate reference projections and local NCI has been following the Force 11 FAIR data principles to make data Findable, Accessible, Interoperable, and Reusable [2]. These principles provide guidelines for a research data repository enable data-intensive science, and enable researchers to answer questions such as how can I improve scientific quality of the data? Is the data usable by my software platform and my tools?

To ensure broader reuse of the data, enable trans-disciplinary integration across multiple domains as well as enabling programmatic access, a dataset must be usable and of value to a broad range of users from different communities [1]. Therefore, a set of standards and 'best practices' for the quality of scientific data products is a critical component in the life cycle of data management

informatics 2017, 4, 45; doi:10.3390/informatics4040045

www.mdpi.com/journal/informatics

<https://doi.org/10.3390/informatics4040045>

DEVELOPMENT PROCESS

TOOLS

PUBLISHING

COMPLIANCE STANDARDS

FILE FORMAT

DATA SERVING: THREDDS, OpenDAP, WMS, etc.
DATA USAGE: Matlab, R, Python, GDAL, etc.

Digital Object Identifiers (DOI) minting,
Making metadata/data available and discoverable online

FILE (GRANULE)-LEVEL	COLLECTION & DATASET-LEVEL
<ul style="list-style-type: none">Climate and Forecasts (CF) ConventionAttribute Convention Dataset Discovery (ACDD)Additional discipline specific standards	Data Management Plans (ISO 19115, ANZLIC, etc.)

Self-describing file formats (e.g., NetCDF, HDF)

7 – Integrar/Publicar

COMBINA DADOS DE DIVERSAS FONTES PARA POSSIBILITAR NOVAS ANÁLISES E INVESTIGAÇÕES.

- O SUCESSO DA INTEGRAÇÃO DE DADOS DEPENDE DO EMPREGO DE BOAS PRÁTICAS DE GESTÃO EM TODO O PROCESSO DE GESTÃO DOS DADOS.
- EXISTEM DIVERSOS CENÁRIOS PARA A INTEGRAÇÃO DE DADOS:
 - INTEGRAÇÃO DE DADOS DE MÚLTIPLOS PROJETOS PARA O TRATAMENTO DE QUESTÕES COMPLEXAS;
 - DADOS ESPARSOS QUE PRECISAM SER COMPLEMENTADOS COM DADOS EXISTENTES PARA POSSIBILITAR A REALIZAÇÃO DE ANÁLISES;

EXEMPLOS:

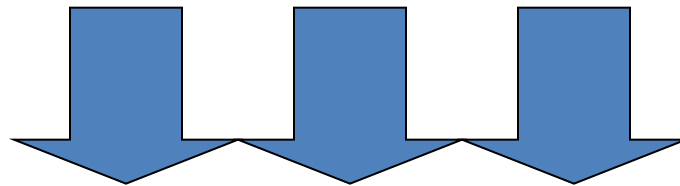
CKAN (DADOS ABERTOS DE GOVERNO – [HTTP://CKAN.ORG](http://ckan.org)) , QUANDL (DADOS FINANCEIROS - [HTTPS://WWW.QUANDL.COM](https://www.quandl.com)), KAGGLE DATASETS. ([HTTPS://WWW.KAGGLE.COM/DATASETS](https://www.kaggle.com/datasets)), MICROSOFT AZURE MARKET PLACE ([HTTP://DATAMARKET.AZURE.COM/BROWSE/DATA](http://datamarket.azure.com/browse/data))

7-Integrar/Publicar

1. ESTRATÉGIAS PARA TORNAR AS INFORMAÇÕES SOBRE OS DADOS DISPONÍVEIS, PARA QUE OS OUTROS POSSAM DESCOBRI-LOS E ACESSÁ-LOS:
 1. AUMENTANDO A VISIBILIDADE DO PROJETO E DE SEUS DADOS:
 - AUMENTO DO POTENCIAL DE USO AMPLO;
 - BENEFÍCIOS PARA A PESQUISA CIENTÍFICA, O APOIO À DECISÃO E A ELABORAÇÃO DE POLÍTICAS PÚBLICAS.
 2. ESTRATÉGIAS PARA GARANTIR O MÁXIMO IMPACTO PARA OS DADOS GERADOS:
 - REGISTRAR O PROJETO EM UM SITE DE DIRETÓRIO DE PROJETOS;
 - DEPOSITAR OS DADOS GERADOS EM UM REPOSITÓRIOS COMPARTILHADOS;
 - ADICIONAR DESCRIÇÕES DOS DADOS (METADADOS) EM SISTEMAS DE ARMAZENAMENTO DISTRIBUÍDO DE METADADOS.

7 - Esforço necessário

- Normalmente, essas atividades envolvem acesso a conjunto de dados independentes, assim:
 - **Usuários devem executar muitas atividades de gerenciamento de dados para extrair, integrar e analisar dados.**
- Necessário uma visão integrada através de todas as redes de dados, permitindo uma busca integrada, análise e visualização através de um conjunto comum de ferramentas e protocolos.



Abordagem de gerenciamento integrado do Ciclo de Vida dos Dados

Conceito 7 – D.O.I. para Dados

- DOI é um acrônimo para "*Digital Object Identifier*", o que significa um "identificador digital de um objeto". Um DOI é um identificador (não um endereço) de uma entidade em redes digitais.
- Um DOI pode ser atribuído a qualquer entidade (objeto) - físico, digital ou abstrato - principalmente para compartilhar em uma comunidade de utilizadores interessados ou na gestão como propriedade intelectual.
- DOI é expresso em URLs (URIs).
- DataCite: organização responsável por implementar D.O.I para dados.

7 – D.O.I. para Dados

Sintaxe

- A sintaxe do DOI especifica a construção de uma cadeia com autoridade e delegação de nomes. Ele fornece um "container" identificador que pode acomodar qualquer identificador existente.
- O DOI tem dois componentes, o prefixo e o sufixo, que juntos formam o DOI, separado pelo caractere "/".

- Não há qualquer limitação sobre o comprimento de um DOI.

- Exemplos:

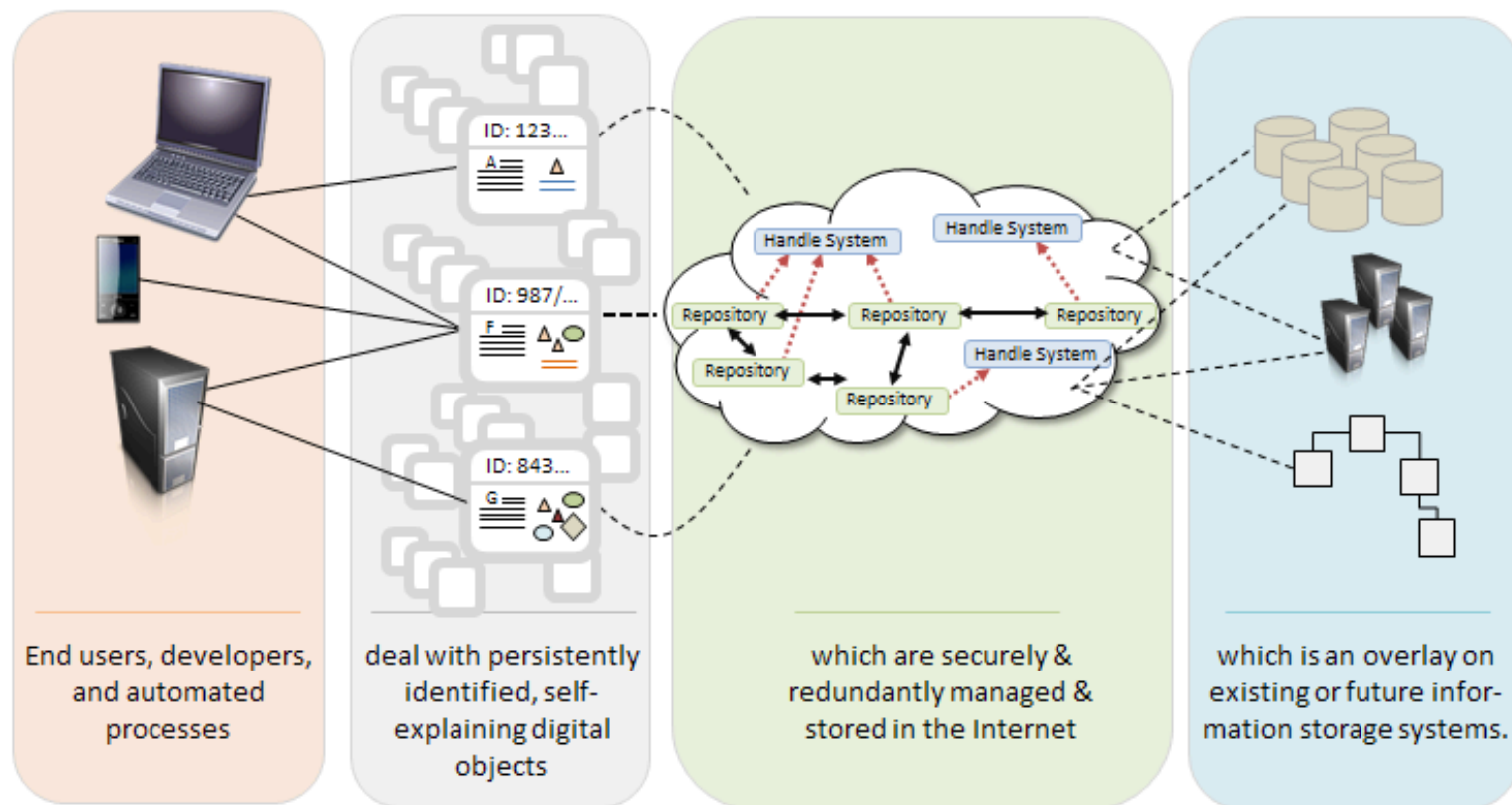
Irino, T; Tada, R (2009): Chemical and mineral compositions of sediments from ODP Site 127-797. Geological Institute, University of Tokyo. <http://dx.doi.org/10.1594/PANGAEA.726855>

Geofon operator (2009): GEFON event gfz2009kciu (NW Balkan Region). GeoForschungsZentrum Potsdam (GFZ). <http://dx.doi.org/10.1594/GFZ.GEOFON.gfz2009kciu>

Denhard, Michael (2009): dphase_mpeps: MicroPEPS LAF-Ensemble run by DWD for the MAP D-PHASE project. World Data Center for Climate http://dx.doi.org/10.1594/WDC/C/dphase_mpeps

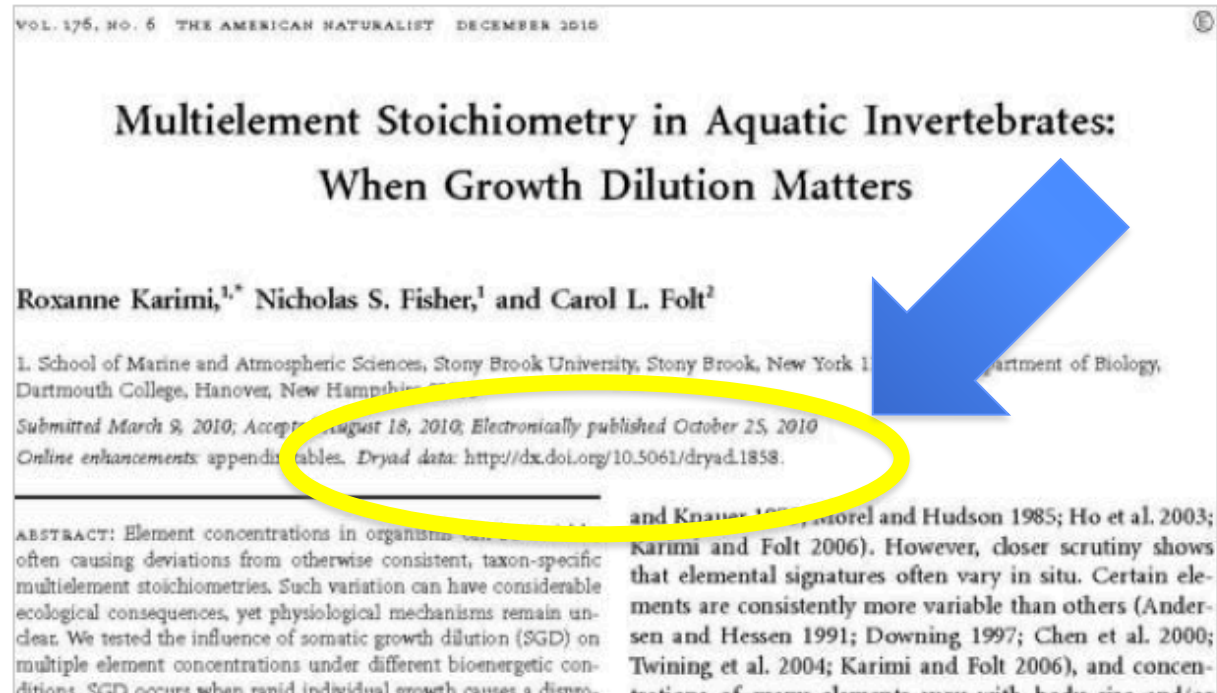
7 – D.O.I. para Dados

Arquitetura



DOI Handbook: Acesso disponível em: https://www.doi.org/doi_handbook/5_Applications.html

7 – D.O.I. para Dados – Citação



When using this data, please cite the original article:

Ally D, Ritland K, Otto SP (2008) Can clone size serve as a proxy for clone age? An exploration using microsatellite divergence in *Populus tremuloides*. *Molecular Ecology* 17(22): 4897-4911.
doi:10.1111/j.1365-294X.2008.03962.x

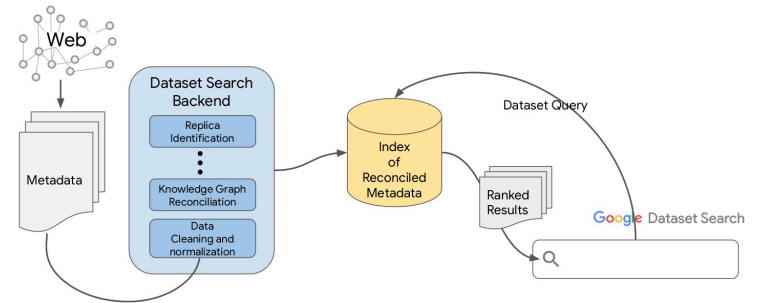
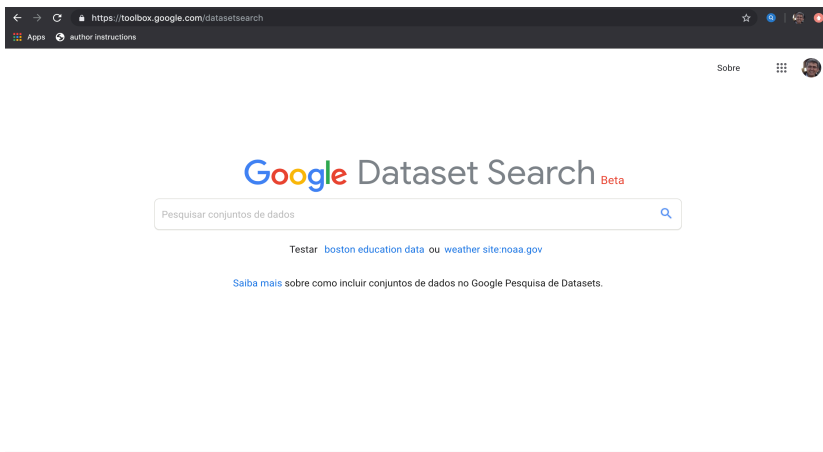
Additionally, please cite the Dryad data package:

Ally D, Ritland K, Otto SP (2008) Data from: Can clone size serve as a proxy for clone age? An exploration using microsatellite divergence in *Populus tremuloides*. Dryad Digital Repository. doi:10.5061/dryad.7898

7 – Exemplo Publicação de Dados

Google Data Search

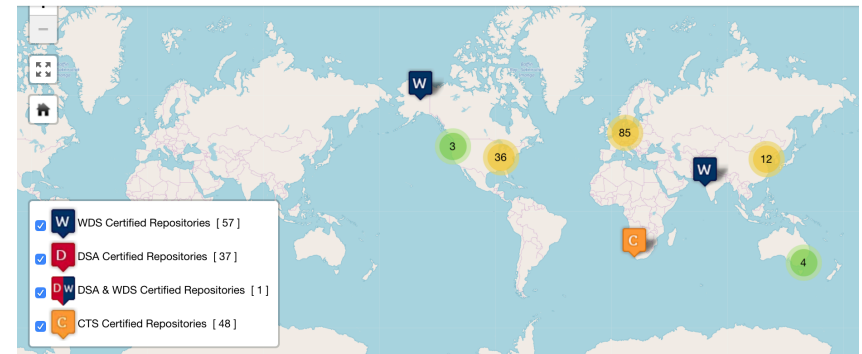
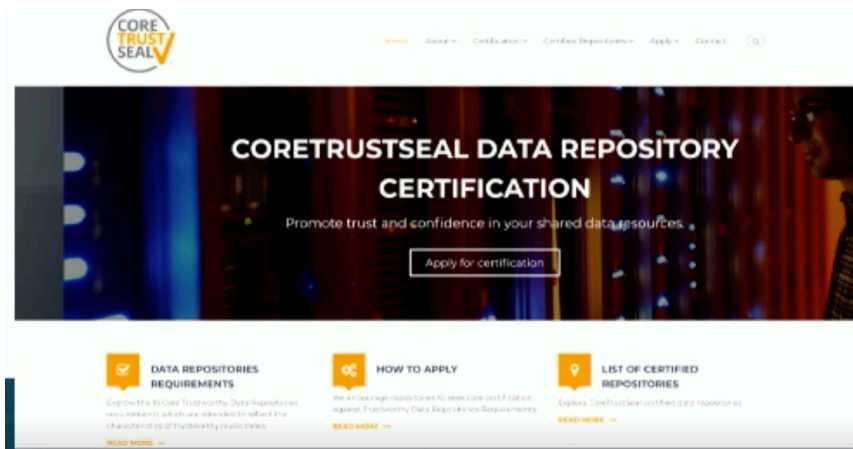
<https://toolbox.google.com/datasetsearch>



Natasha Noy and Matthew Burgess and Dan Brickley , Google Dataset Search: Building a search engine for datasets in an open Web ecosystem, WebConf 2019. available at: <https://ai.google/research/pubs/pub47845>

Confiabilidade para reuso dos dados

- Se alguém entregar seus dados científicos, o que é necessário para convencê-lo que os dados estão corretos ?
- Se requer um sistema complexo para executá-lo e que não tem acesso, o que precisa para confiar dos dados ? Você conhece quais são as suposições e dependências existentes ?
- Quanto você poderá confiar que os mesmo dado estará disponível por um longo período de tempo ?



BOAS PRÁTICAS NA GESTÃO DE DADOS

Integrar



Analisar

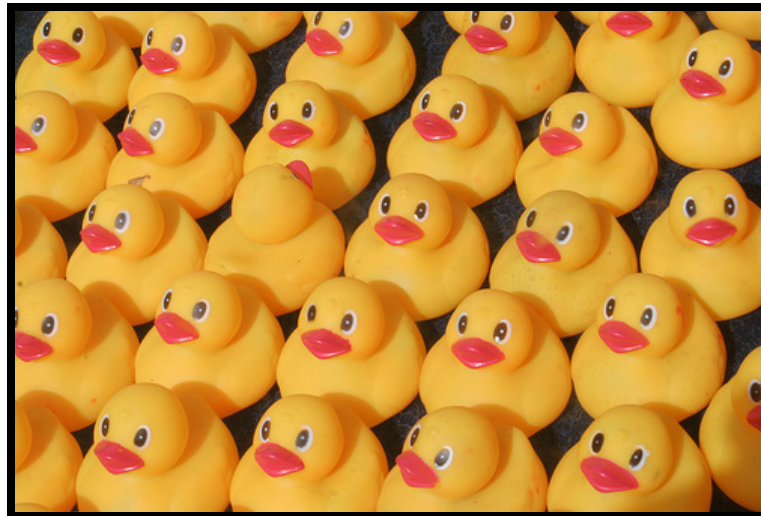


Planejar

8 - Análise

b) Reprodutibilidade

- **A reprodutibilidade é a chave para os métodos científicos;**
- **Processos complexos são mais difíceis de serem reproduzidos;**
- **Boa documentação é essencial para a reprodutibilidade;**



8 - Análise

FLUXOS DE TRABALHO (*WORKFLOW*)

- A FORMALIZAÇÃO DOS METADADOS DO PROCESSO;
- DESCRIÇÃO PRECISA DO PROCEDIMENTO CIENTÍFICO;
- TRÊS COMPONENTES:
 - INPUTS: INFORMAÇÃO E/OU MATERIAL NECESSÁRIO;
 - OUTPUTS: INFORMAÇÃO OU MATERIAL PRODUZIDO E POTENCIALMENTE UTILIZADO COMO INPUT EM OUTROS PASSOS;
 - REGRAS DE TRANSFORMAÇÃO/ALGORITMOS.

8 - Análise

BOAS PRÁTICAS

- OS CIENTISTAS DEVEM DOCUMENTAR OS FLUXOS DE TRABALHO USADOS NA CRIAÇÃO DE RESULTADOS:
 - PROVENIÊNCIA DOS DADOS;
 - ANÁLISES E PARÂMETROS UTILIZADOS;
 - CONEXÕES ENTRE ANÁLISES POR MEIO DOS INPUTS (ENTRADAS E OUTPUTS (SAÍDAS)).
- A DOCUMENTAÇÃO PODE SER INFORMAL (EX: FLOWCHARTS, *COMMENTED SCRIPTS*) OU FORMAL (EX: KEPLER, VISTRAILS).

8 - Análise

- Workflow Formal

VisTrails

VisTrails VCDAT

History Tree

VisTrails Shell

```
VisTrails shell running Python 2.6.4 (r264:75821M, Oct 27 2009, 19:48:32)
[[GCC 4.0.1 (Apple Inc. build 5493)] on darwin.
Type "copyright", "credits" or "license" for more information on Python.
>>> import vcs, cdm2
>>> cdat = load_package('DAT')
>>> cdmsfile = cdm2.open('/home/amanuele/src/cdat_bin/sample_data/clt.nc')
>>> data = cdm2file['clt']
>>> q = cdat.quickplot()
>>> q.dataset = data
>>> run()
```

Visualization SpreadSheet

Drag & drop
components
from this list

www.vistrails.org

Kepler Software

Kepler Software

Workflow

CT Director

Timed Plotter

XY Plotter

Integrator

Integrator

$r = 2$
 $a = 0.1$
 $b = 0.1$
 $d = 0.1$

$\frac{dn_1}{dt} = r \cdot n_1 - a \cdot n_1 \cdot n_2$

$\frac{dn_2}{dt} = -d \cdot n_2 + b \cdot n_1 \cdot n_2$

This model shows the solution to the classic Lotka-Volterra predator-prey dynamics model. It uses the Continuous Time domain to solve two coupled differential equations, one that models the predator population and one that models the prey population. The results are plotted as they are calculated showing both population change and a phase diagram of the dynamics.

Rich Williams, 2003, NCEAS

Actors in
workflow

kepler-project.org

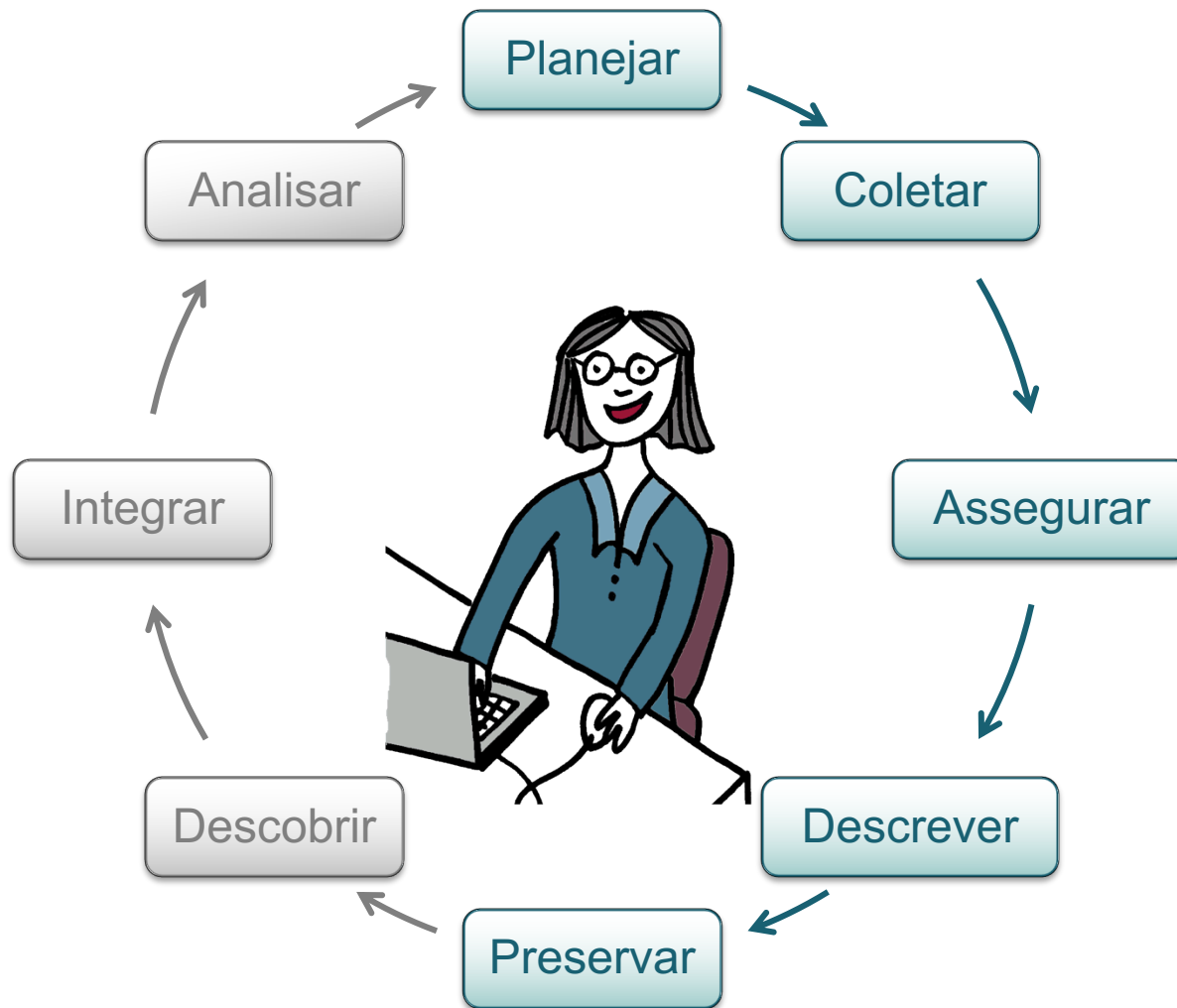
Software - Próxima fronteira



- Similar ao processo de publicação dos Dados
- O software que Analisa os dados precisa ser preservado para efeitos e reprodutibilidade
- Considerações: <https://www.usgs.gov/media/files/usgs-software-planning-checklist>
 - Categorias de disponibilização – Informal e formal
 - Formal: Disclaimers (Provisório/Aprovado)
 - Licenses – Código pode ser de domínio público e/ou incluir restrições de terceiros;
 - Estratégias de Documentação - <https://github.com/usgs/best-practices>
 - Code Reviews (PII/Security, scientific verification, standards)- <https://github.com/usgs/best-practices>
 - Obtenção de um DOI <https://github.com/usgs/best-practices/blob/master/doi.md>



Gestão de Dados



Exercício - Desafio

1- Considere um domínio de aplicação envolvendo dados de seu interesse. Defina um plano de gestão de dados considerando a ferramenta dmp.cdlib.org.

2- Para esse domínio selecione um arquivo em formato csv, com pelo menos 2 atributos. Defina os metadados para esse arquivo utilizando as propriedades recomendadas do Google Search.

Referências:

- Diretrizes para DataSets <https://developers.google.com/search/docs/data-types/dataset>

- Natasha Noy and Matthew Burgess and Dan Brickley , Google Dataset Search: Building a search engine for datasets in an open Web ecosystem, WebConf 2019. available at: <https://ai.google/research/pubs/pub47845>

Conceitos de Big Data

PCS5787 – Ciência dos Dados
Prof. Dr. Pedro Luiz Pizzigatti Corrêa
24 de Setembro de 2020