

Dilúvio de Dados

Uma Boia à Vista,
ou seria um Cruzeiro?

MIRELLA M. MORO
UFMG



**SBBD
2023**

Quantidade



Comunidade

Mar 2012
SBC



Carol

Quantidade



Comunidade



Karin

**“Vida de
cientista
significa estar
sempre
aprendendo”**





Data deluge

Google

data deluge



Images

Example

Pdf

Wikipedia

Synonym

Videos

News

Zoho

Shopping

About 14,100,000 results (0.55 seconds)

Definition. The data deluge refers to the situation where the sheer volume of new data being generated is overwhelming the capacity of institutions to manage it and researchers to make use of it.

Text

Images

Documents

Websites

Detect language

English

Italian

F



Portuguese

English

Spanish



deluge



'delyoo(d)ZH



6 / 5,000



dilúvio



Definitions of deluge

Noun

1

a severe flood.

"this may be the worst deluge in living

Translations of deluge

Noun

Frequency ?

o dilúvio

flood, deluge, hailstorm



a inundação

flood, flooding





Dilúvio de Dados

Dilúvio de Dados

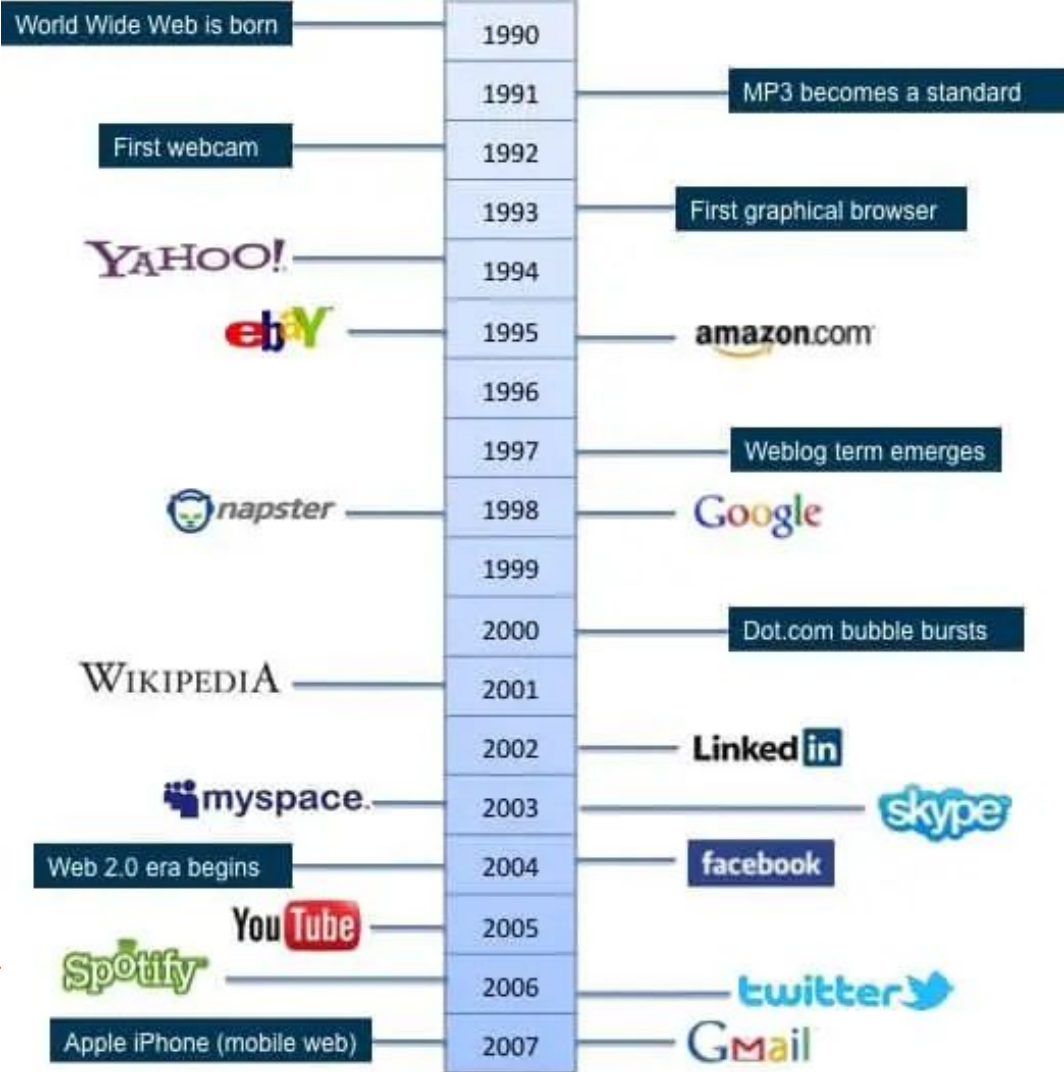
Uma Boia à Vista,
ou seria um Cruzeiro?



1. Uma Boia à Vista

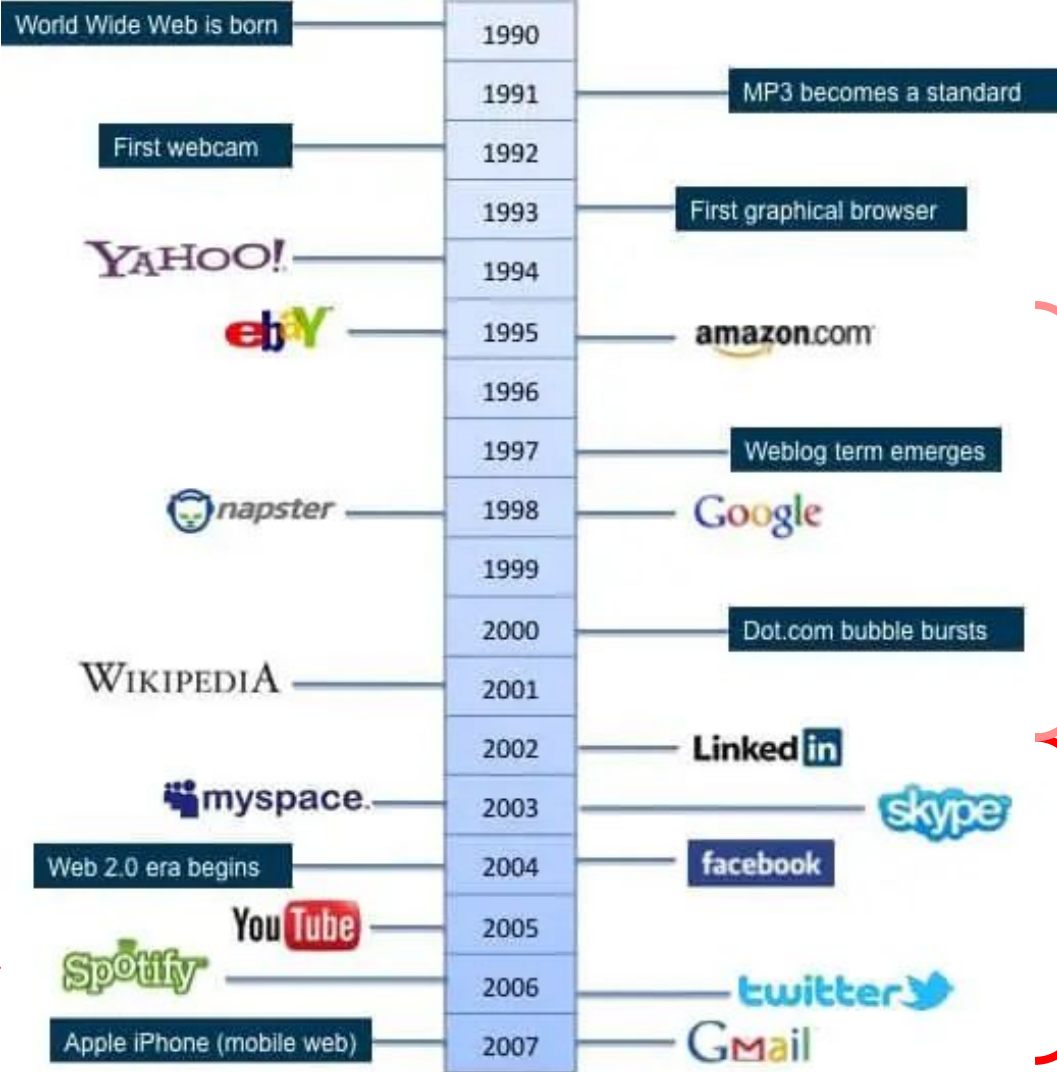


*“Entender de tecnologia é
entender de história e evolução”*



Fonte:

<https://harbott.com/visual-timeline-of-internet-milestones>



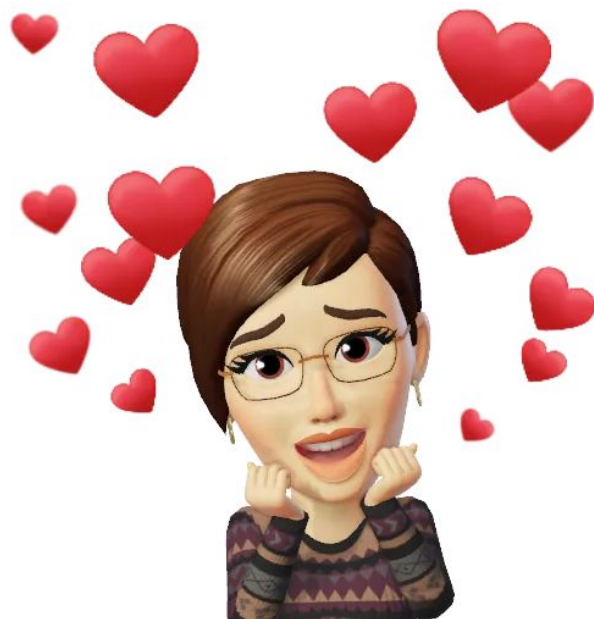
B+M
@UFRGS
PhD
@UCR



Criar + Consumir Conteúdo Web

- Cientistas
- Profissionais
- Pessoa comum

Dilúvio de Dados





XML

Extensible Markup Language

Introduction



EM 2007,
BONS VENTOS TE LEVEM
A JOÃO PESSOA!

SBBD SBES 2007

Tutorial

- XML data everywhere
 - Web services (RSS feed format)
 - Integration (e-commerce, stocks)
 - Open standards for file exchange (FpML, ACORD)
- XML database
 - XML native engine
 - Specialized storage and retrieval
 - XML-enabled Relational DBMS
 - Shredding or text

XML Databases

XML Document = sequence of elements that enclose text values and other elements

```
<book >
  <title>"Adventures of Huckleberry Finn"</title>
  <author>"Mark Twain"</author>
  <year>2002</year>
  <otherInfo>
    <isbn>0142437174</isbn>
    <collection>"Penguin Classics"</collection>
  </otherInfo>
</book>
```

XML Databases

XML Document = sequence of elements that enclose text values and other elements

Element name between < >



Text values between elements



```
<book >  
  <title>“Adventures of Huckleberry Finn”</title>  
  <author>“Mark Twain”</author>  
  <year>2002</year>  
  <otherInfo>  
    <isbn>0142437174</isbn>  
    <collection>“Penguin Classics”</collection>  
  </otherInfo>  
</book>
```

**“Menos é mais
na vida e no
slide”**



XML Nativo

Minicurso

Mirella M. Moro

Departamento de Ciência da Computação

Universidade Federal de Minas Gerais

mirella@dcc.ufmg.br



XXIII SIMPÓSIO BRASILEIRO DE BANCO DE DADOS

XXII SIMPÓSIO BRASILEIRO DE ENGENHARIA DE SOFTWARE

Campinas, Brasil - 13 a 17 de Outubro

XML Nativo

- XML

A linguagem que apareceu para resolver todos os problemas do mundo

- É... tipo assim...

XML está por toda a parte...

- Advertising: [adXML](#) place an ad onto an ad network or to a single vendor
- Banking: [MBA](#) Mortgage Bankers Association of America → credit report, loan file, underwriting...
- Directories: [dirXML](#) Novell's Directory Services Markup Language
- **Literature:** [Gutenberg](#) convert the world's great literature into XML
- Geospatial: [ANZMETA](#) distributed national directory for land information
- **Healthcare:** [HL7](#) DTDs for prescriptions, policies & procedures, clinical trials
- Human Resources: [XML-HR](#) standardization of HR/electronic recruiting XML definitions
- International Dvt: [IDML](#) improve the mgt. and exchange of info. for sustainable development
- Math: [MathML](#) Mathematical Markup Language
- **News:** [NewsML](#) creation, transfer and delivery of news
- Surveys: [DDI](#) Data Documentation Initiative, "codebooks" in the social and behavioral sciences
- **Travel:** [openTravel](#) information for airlines, hotels, and car rental places
- Voice: [VoxML](#) markup language for voice applications
- Wireless: [WAP](#) Wireless Application Protocol, wireless devices on the World Wide Web
- Weather: [OMF](#) Weather Observation Markup Format
- Web Servers: [apacheXML](#) parsers, XSL, web publishing
- ...

DESMISTIFICANDO XML: DA PESQUISA À PRÁTICA INDUSTRIAL

Mirella Moura Moro – UFMG
Vanessa Braganholo – UFRJ



@SBBD 2007

Bento Gonçalves – RS – Brasil

—
“Tutorial +
Minicurso +
JAI =
trifecta de
recém doutor/a”



Solução ou Problema?

122

- DBLP
- XML \rightarrow +5500
- XML query \rightarrow 604

MORO, BRAGANHOLO,
DORNELES, DUARTE,
GALANTE, MELLO.

XML: Some Papers in a Haystack.
SIGMOD Record, 2009.

Renata, Vanessa,
Ronaldo, Carina



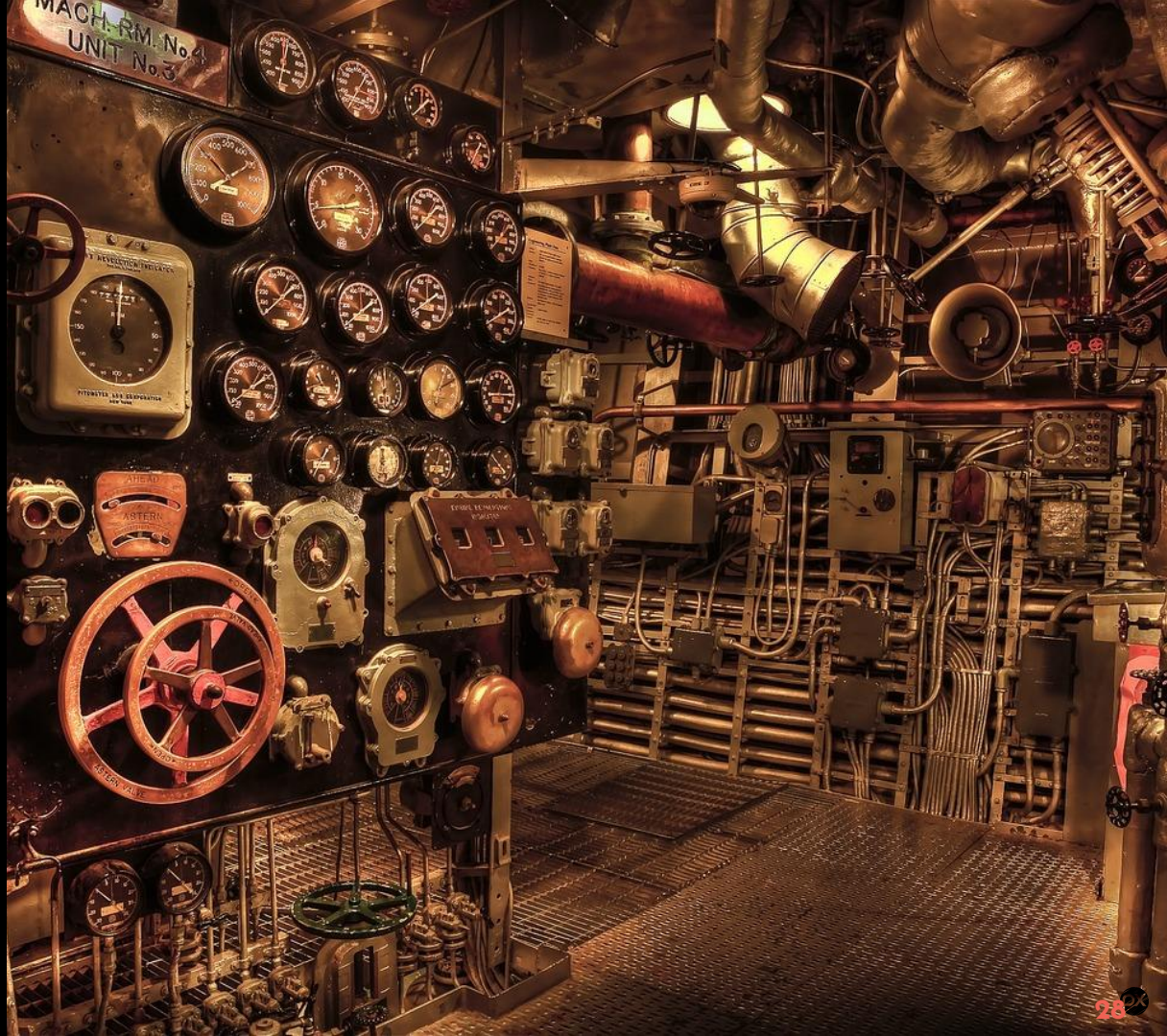
Dilúvio de Dados

Uma Boia à Vista

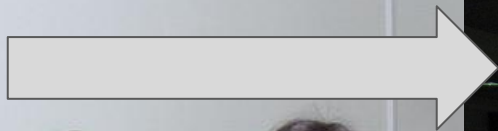
Integração de dados
Padrões
Mídias diferentes



Engenharia de Dados



Cientista da
Computação



Engenheira
de Dados



Setembro 2002

Abril 2007

Recém Doutora

Transição

Recomendação (*streams*)

Busca (*keyword search*)

NoSQL



Palazzo



Alberto



Clodoveu



Altigran

**“Acolhimento
das gerações
mais novas faz
a diferença!”**



2. Ou seria um Cruzeiro





Dilúvio de Dados

Social



Social Professional Networks

 Redes Sociais Acadêmicas

**“Ciência até
surge do nada.
Mas só se
fortalece com
colaboração.”**





The Brazilian Portal of Science and Technology

Alberto H. F. Laender, Mirella M. Moro, Altigran S. Silva, Clodoveu A. Davis Jr., Marcos André Gonçalves, Renata Galante, Allan J. C. Silva, Carolina A. S. Bigonha, Daniel Hasan Dalip, Eduardo M. Barbosa, Eduardo N. Borges, Eli Cortez, Peterson Procópio Jr., Rafael Odon de Alencar, Thiago N. C. Cardoso, Thiago Salles

SEMISH 2011

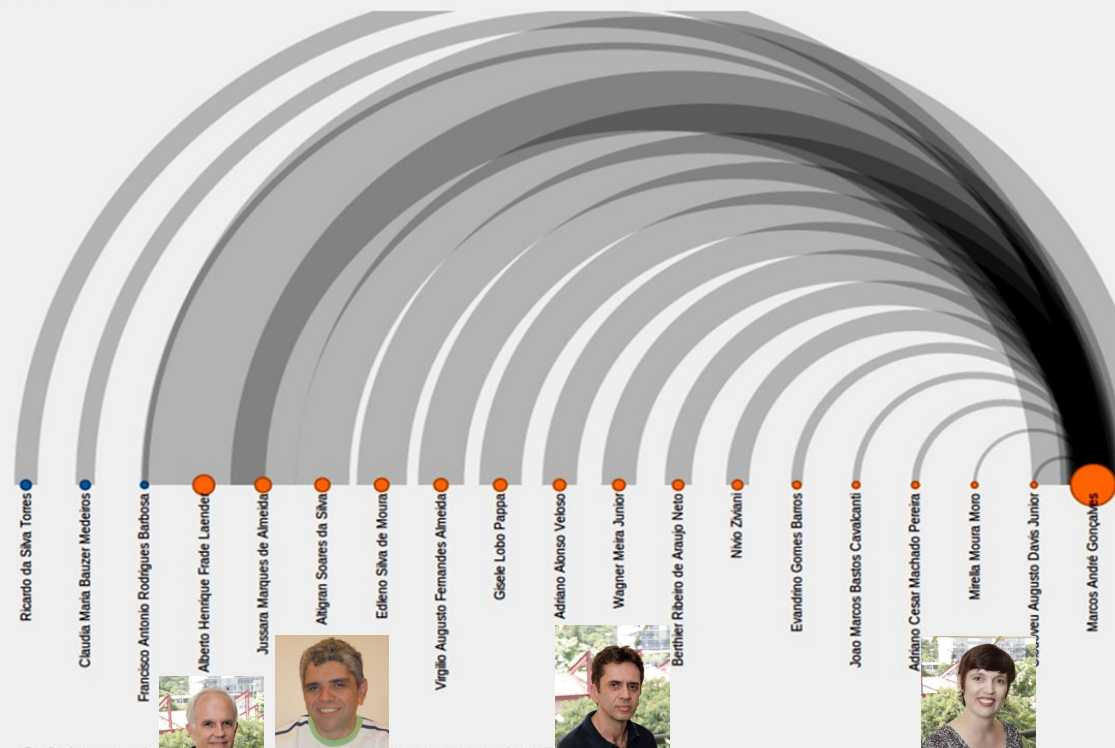


winweb

Instituto Nacional de Ciência e Tecnologia para a Web

Grafo de Coautoria de um Pesquisador

Grafo de Co-autoria



Grafo de co-autoria de um pesquisador, considerando apenas membros deste INCT. Os nós representam pesquisadores do INCT e os arcos as relações de co-autoria. A espessura do grafo é dada pelo número de publicações em conjunto.

- Pesquisadores do mesmo INCT
- Pesquisadores de outros INCTs



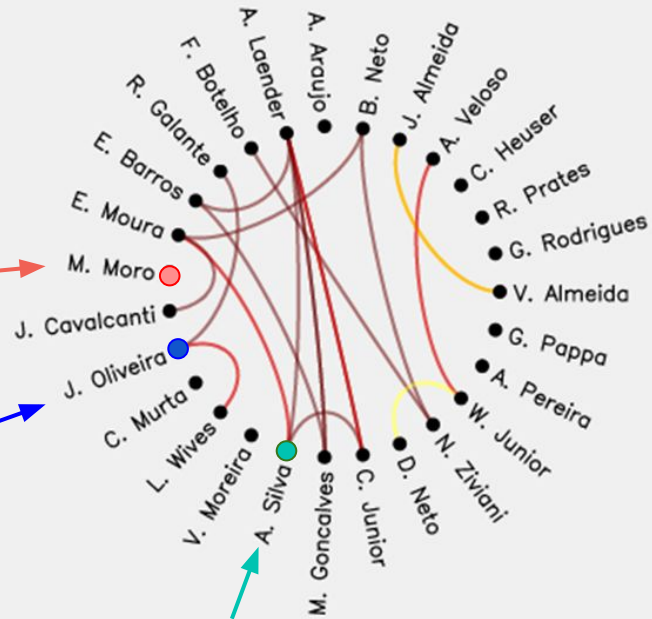
CiênciaBrasil: Groups - período

2005

Filtrar conexões entre

20

Ok

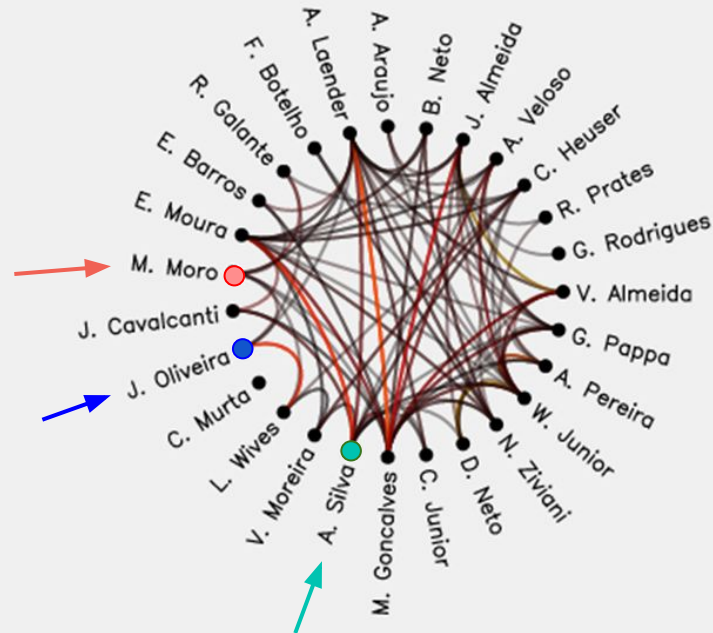


2005:2010

Filtrar conexões entre

20

Ok



Arquitetura

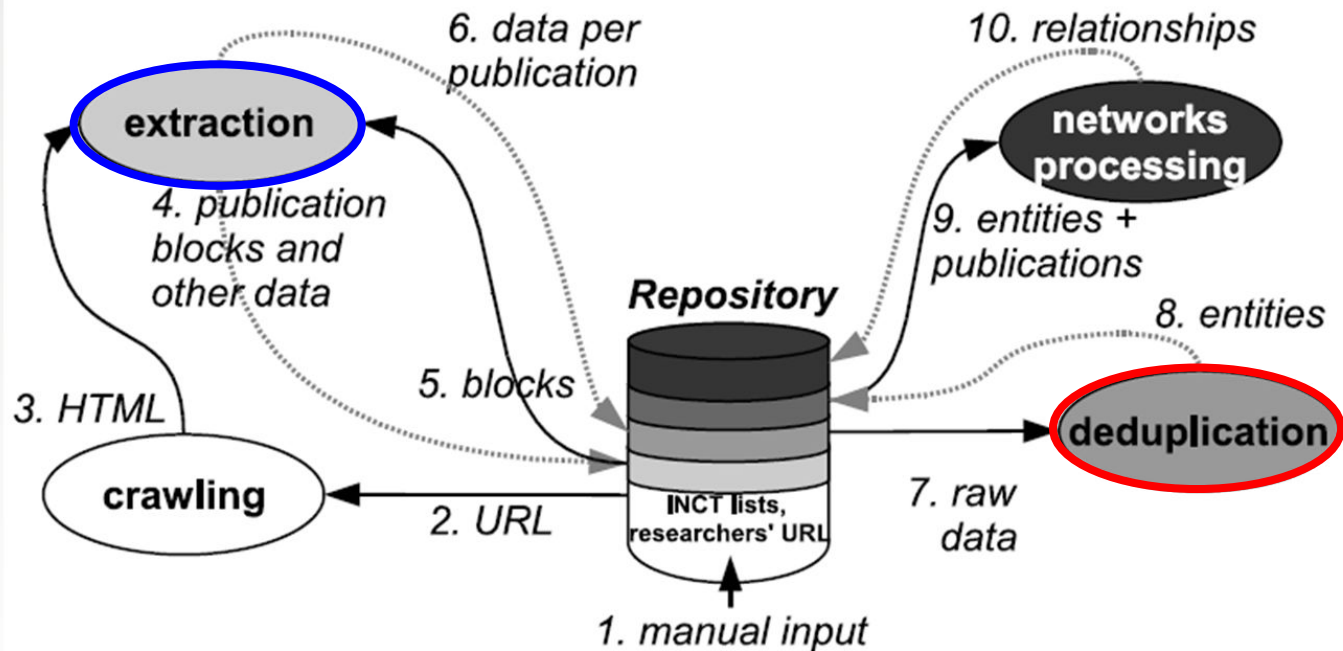
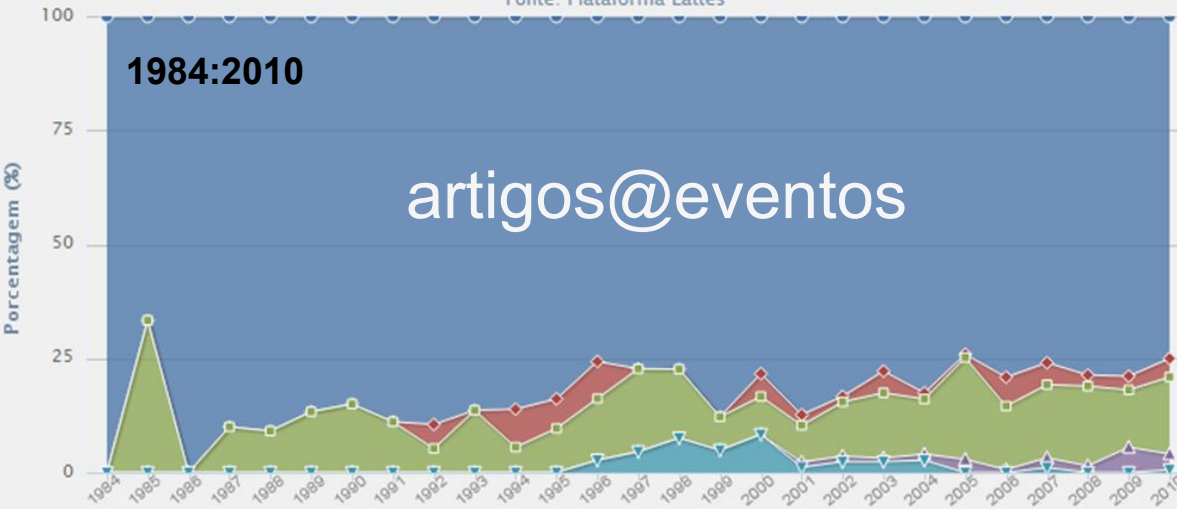
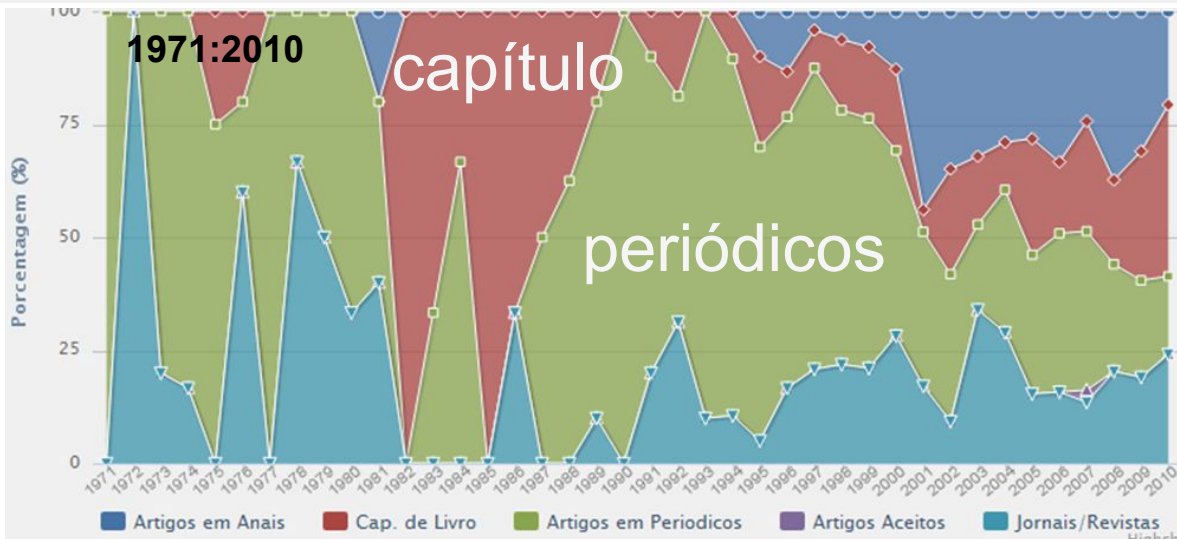


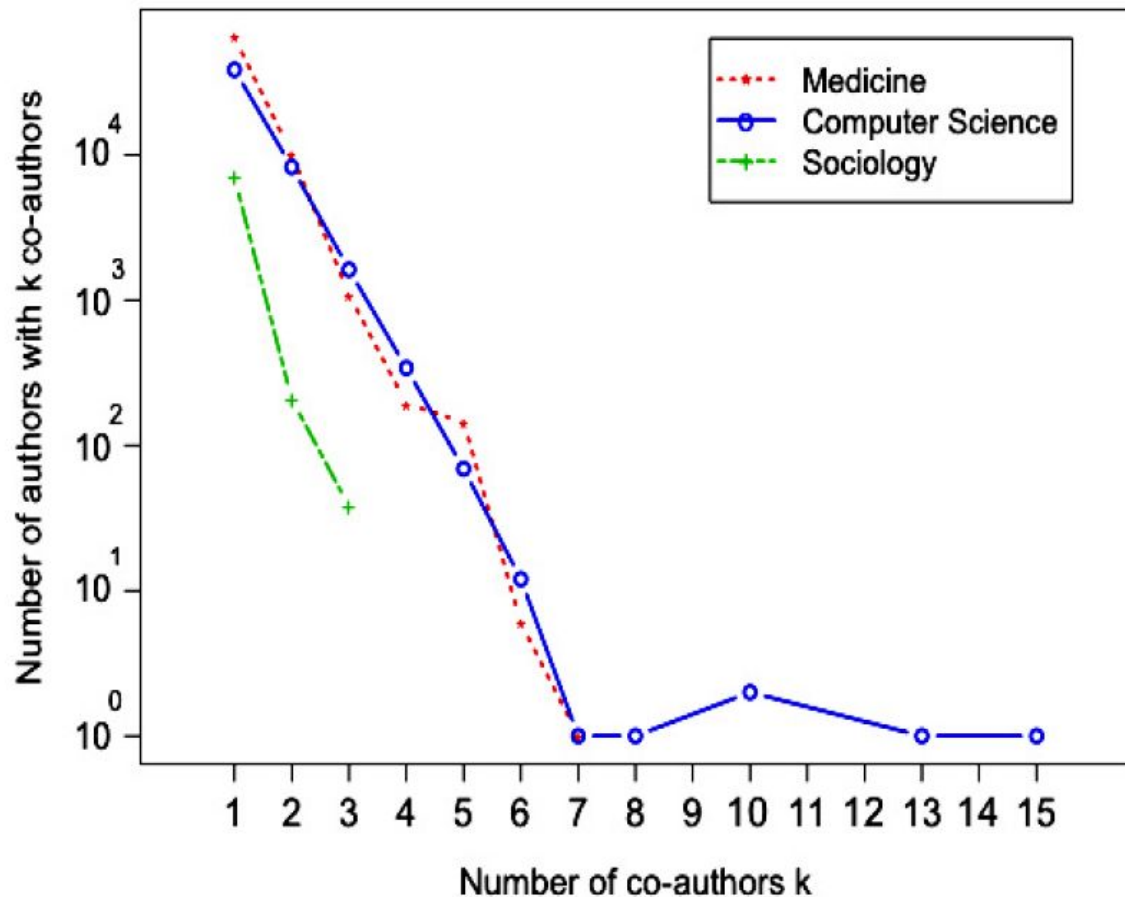
Figure 1. *CiênciaBrasil* Architecture



INCT
Engenharia de
Software



INCT
Ciências Sociais



Distribution of numbers of co-authors for researchers in each area

In **Sociology**: from 7,195 publications, 84% have only one author

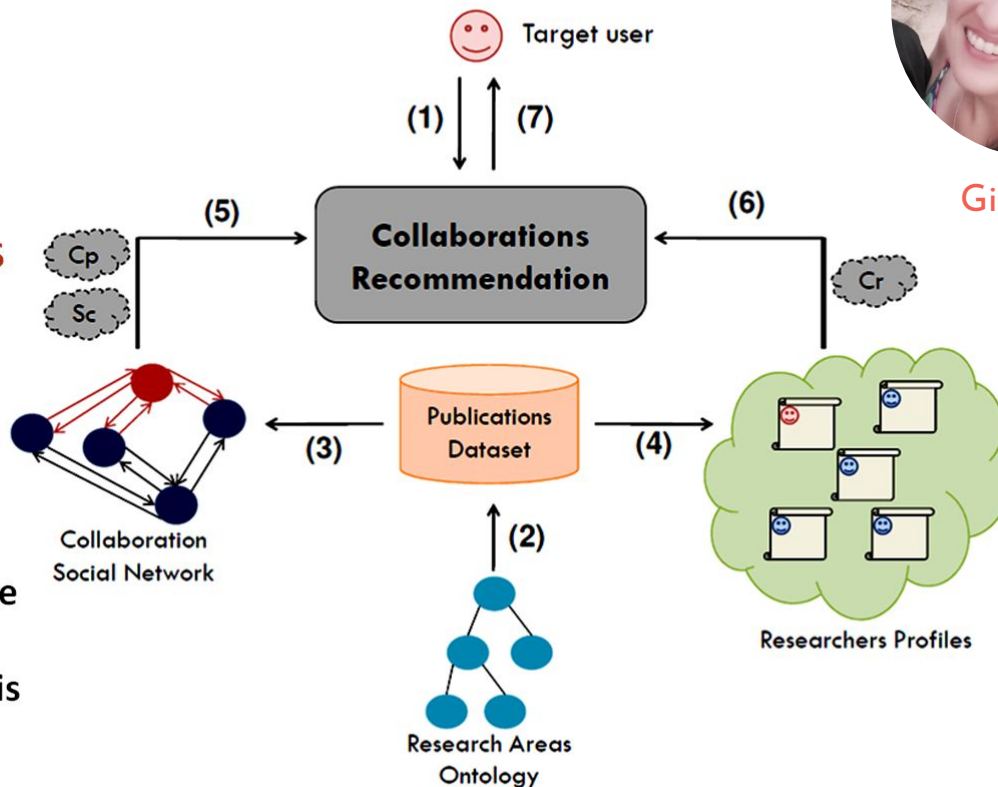
**“Analise os dados
antes de sair
criando redes
complexas em
cima :-)”**



Recomendação de Pessoas

- **CORALS**: a **C**ollaboration **R**ecommender for **A**cademic **s**ocial **n**etwork**S**

- **ReDisseR**: **R**ecomendação e **D**isseminação em **R**edes **S**ociais

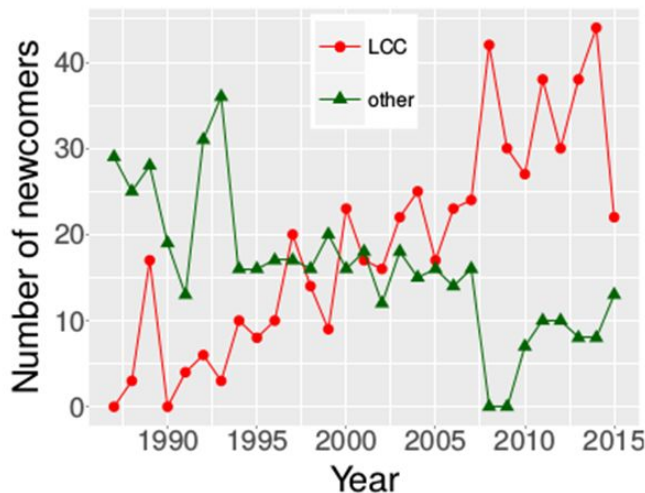


Giseli Lopes



Michele Brandão

Collaborations and Newcomers



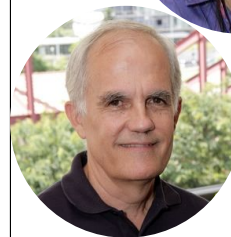
The Collaboration Network of the
Brazilian Symposium on Databases
30 Editions of History

JBCS 23:10 (2017)

10.1186/s13173-017-0059-6

SBBD initially received several authors with no link
to its Largest Connected Component

It has become more common to join the
symposium by collaborating with someone
already connected to it





Social Professional Networks

Coautoria foi só o início

Intensidades Diferentes

Forças dos
relacionamentos

Pessoas influentes

Mobilidade

Ranking





PPG

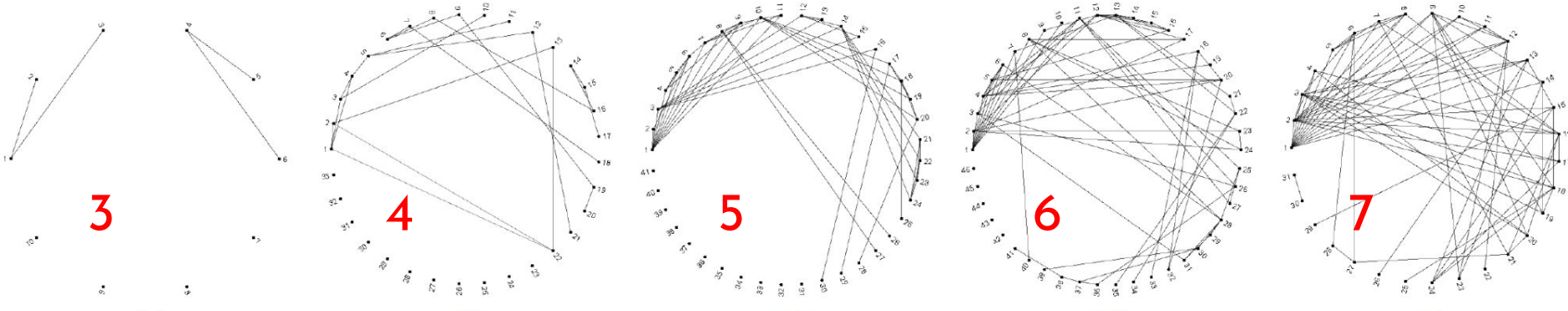


Software@GitHub

Ranking PPG



Colaborações internas ao PPG



Ranking Strategy for Graduate Programs Evaluation
@ ICITA 2011



Aspectos Temporais para Medir a Força da Colaboração no GitHub

Natércia A. Batista, Michele A. Brandão,
Ana Paula C. da Silva, Mirella M. Moro



Best Short Paper

SBBD2017

BRACIS2017

KDMiLE

ENIAC

STIL

Utilização de Redes Heterogêneas para Medir a Força dos Relacionamentos no GitHub



Gabriel P. Oliveira, Natércia A. Batista, Michele A. Brandão, Mirella M. Moro
Universidade Federal de Minas Gerais (UFMG)



Social professional networks: A survey and taxonomy

Michele A. Brandão  , Mirella M. Moro 

Show more 

 Outline |  Add to Mendeley  Share  Cite

<https://doi.org/10.1016/j.comcom.2016.12.011>

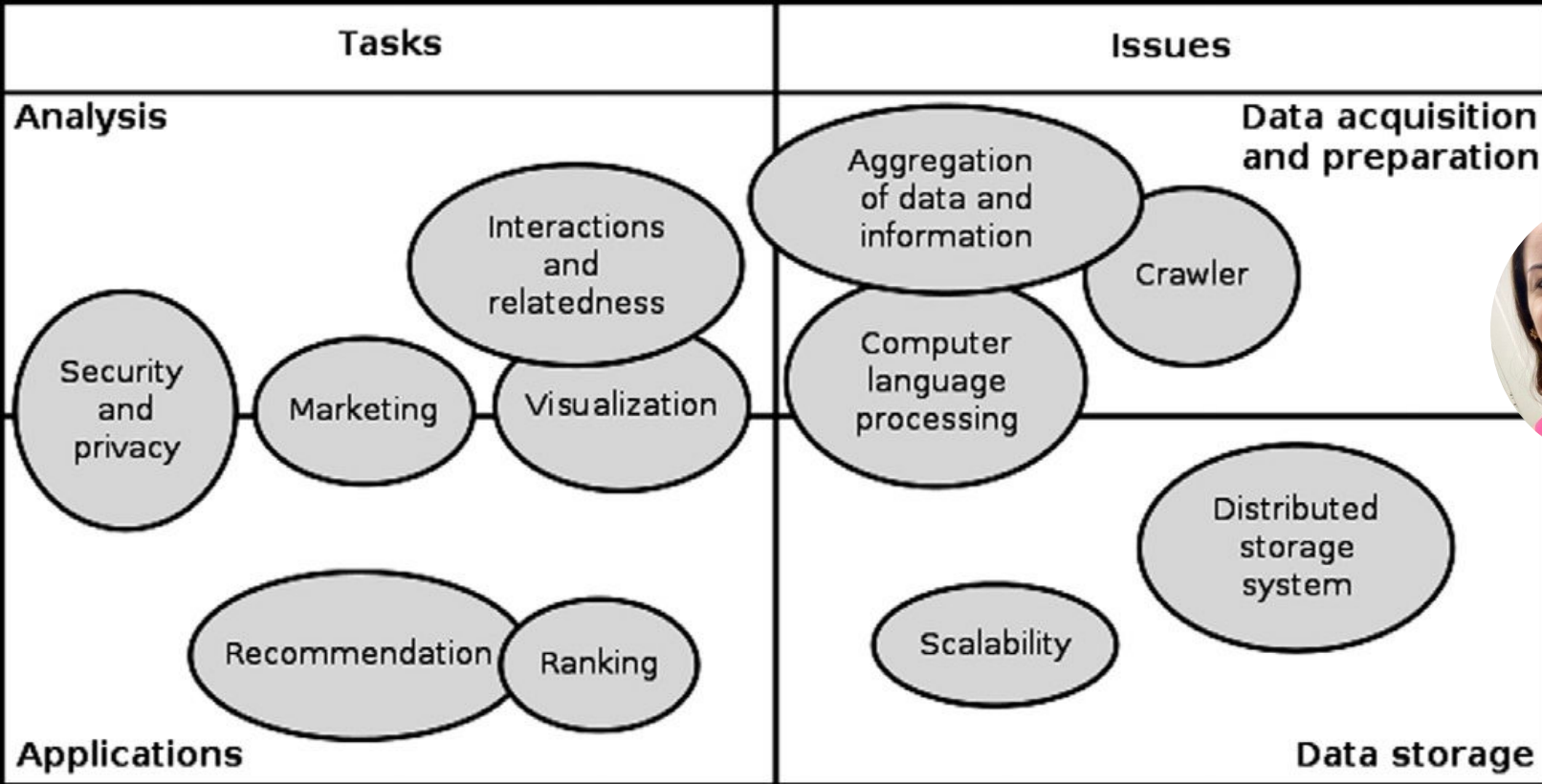


Fig. 5. Main social networks topics, in which **tasks** refer to **using** social networks to solve problems, and **issues** address problems related to **managing** social networks.

**“Tese de
doutorado tem
de produzir
survey em
periódico”**





Social Professional Networks

Coautoria + PPG
Software

**“Às vezes a
oportunidade
realmente bate a
sua porta”**



Ouvir

Quem você orienta
2017



Laís Mota (M)

Ouvir

Quem você orienta
2017



Laís Mota (M)



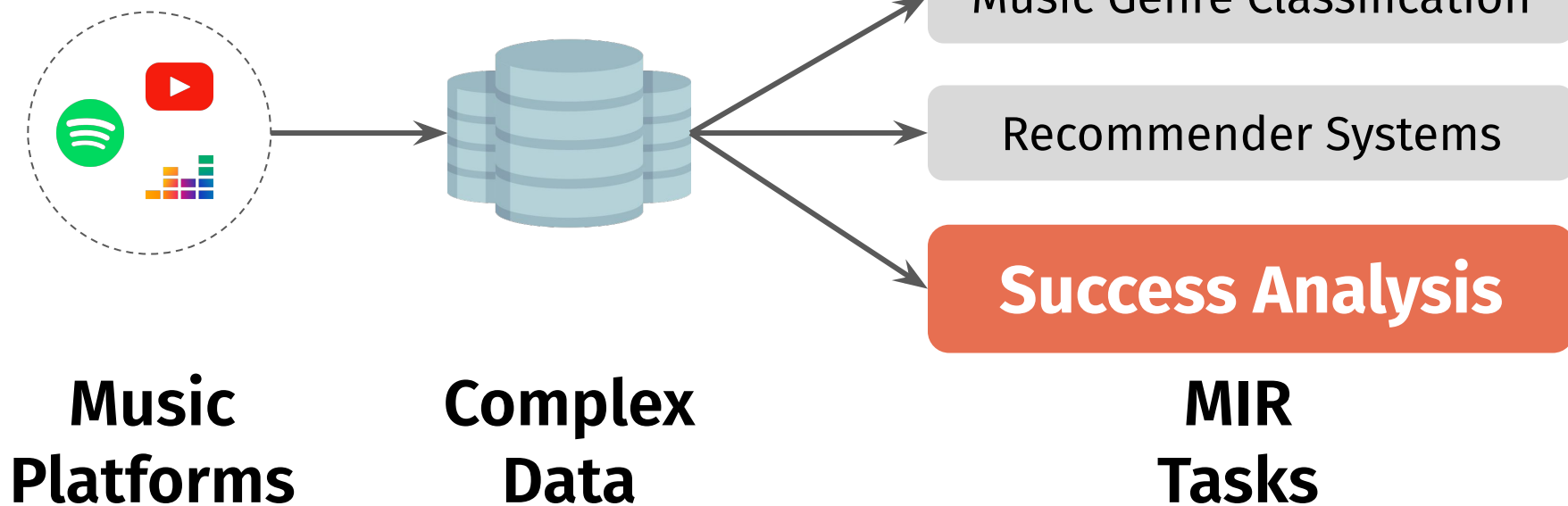
Mariana O S Silva (M)

The background of the slide features a close-up, slightly blurred view of a piano keyboard and an open sheet of music. The sheet music is on aged, cream-colored paper and contains several staves with musical notation, including treble and bass clefs, notes, and rests. The piano keys are visible in the lower portion of the frame, with the white keys being more prominent than the black ones. The overall lighting is soft and warm, creating a professional and artistic atmosphere.

Ciência de Dados na indústria da Música



Music-related Data



Sucesso





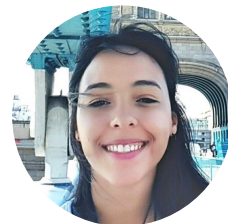
Sucesso
depende
dos dados



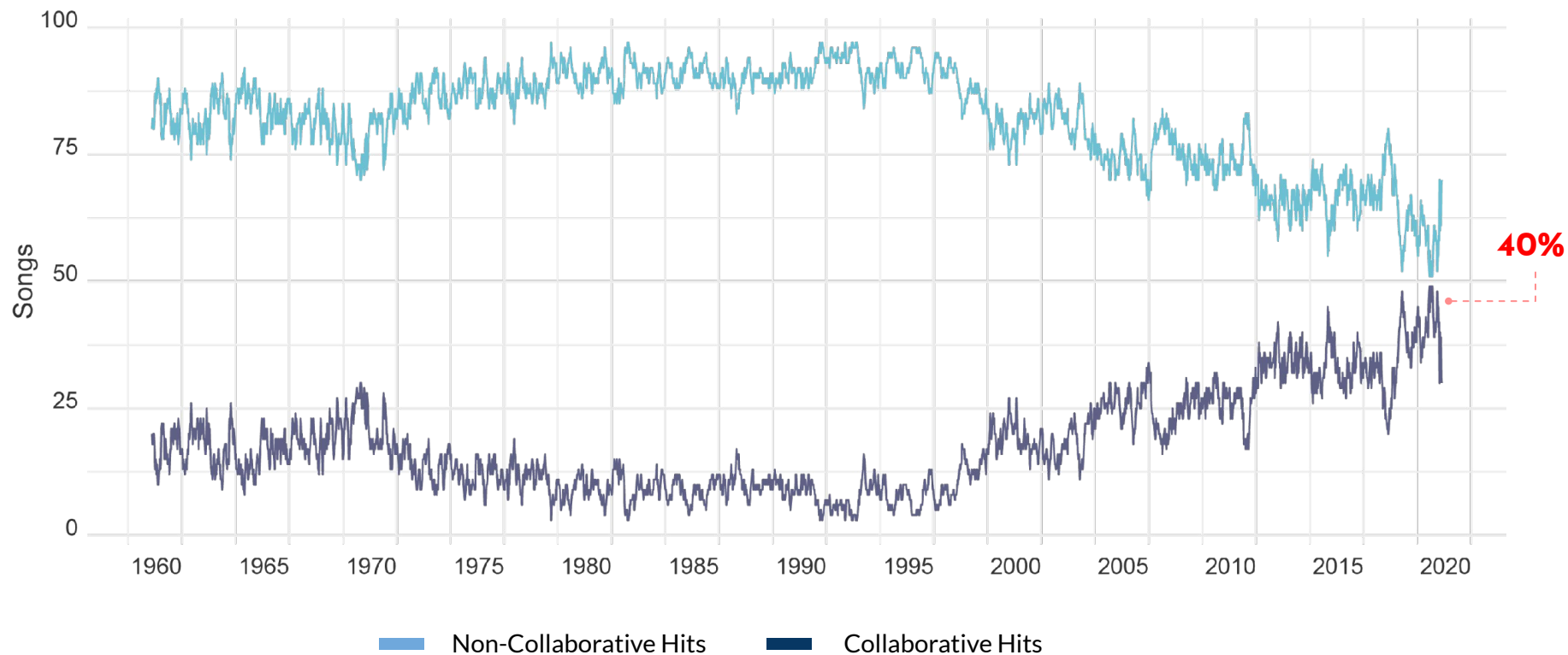
Colaboração



Collaborative Hits



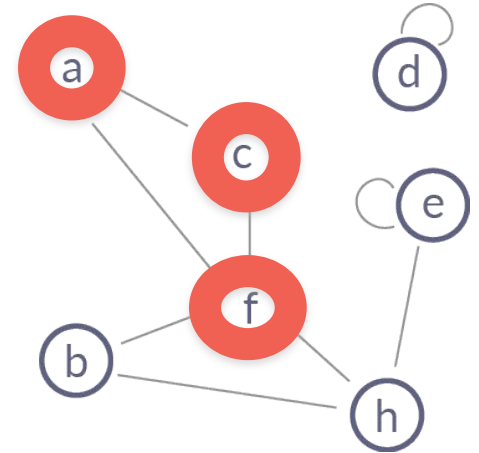
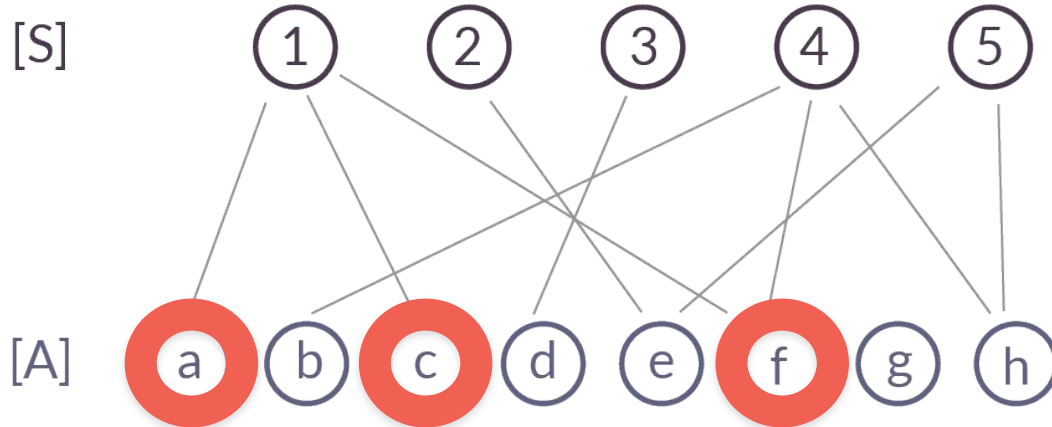
Billboard Hot 100 songs (1958 – 2018)



Social Network Modeling



SINGLES

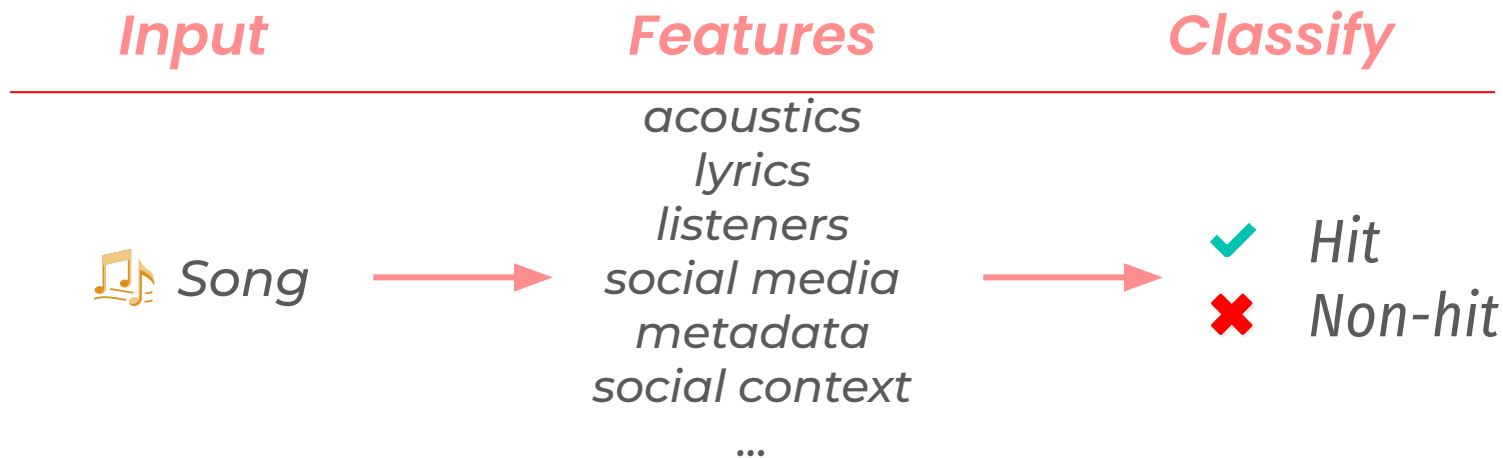


ARTISTS



HIT SONG PREDICTION

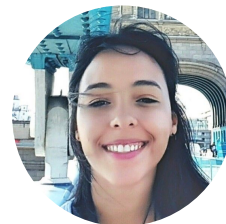
Binary Classification



Hit = Billboard Hot 100

Non-hit = artists' songs \notin BH100

MULTI-PERSPECTIVE **FEATURES**



Song

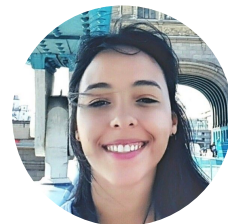


Album



Artist

COLLABORATION-BASED **FEATURE**

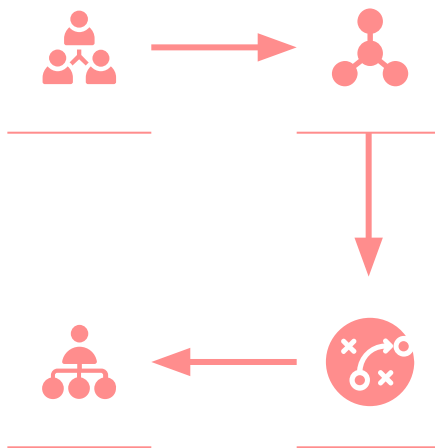


Collaboration Network

Weighted edges
connect artists who
have collaborated

Collaboration Profiles

Well-defined
artists groups



Topological Metrics

Compute seven
topological metrics

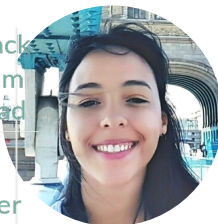
Clustering (K-Means)

Clustering similar
artists, based on
collaboration
patterns

■ Diverse

▲ Regular

● Absent



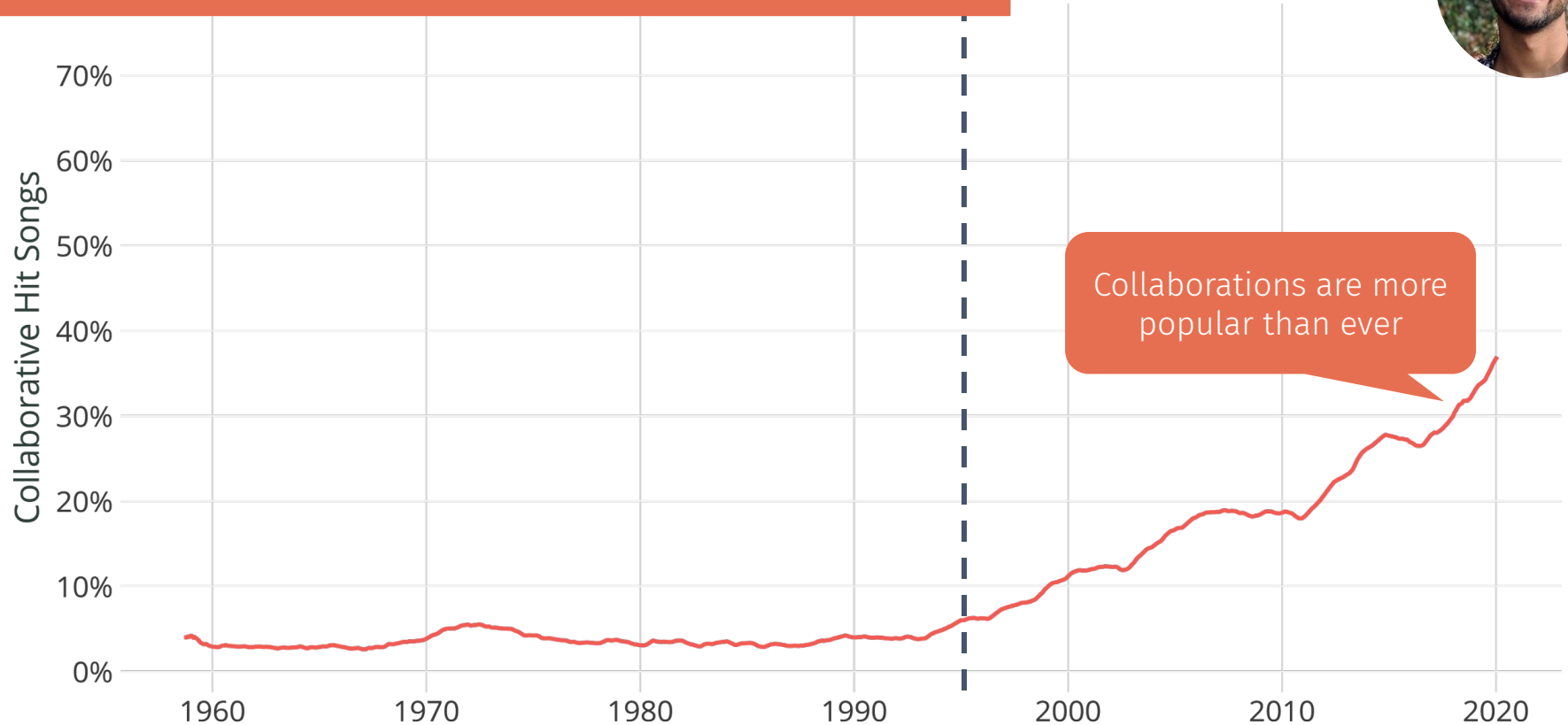
■ Diverse

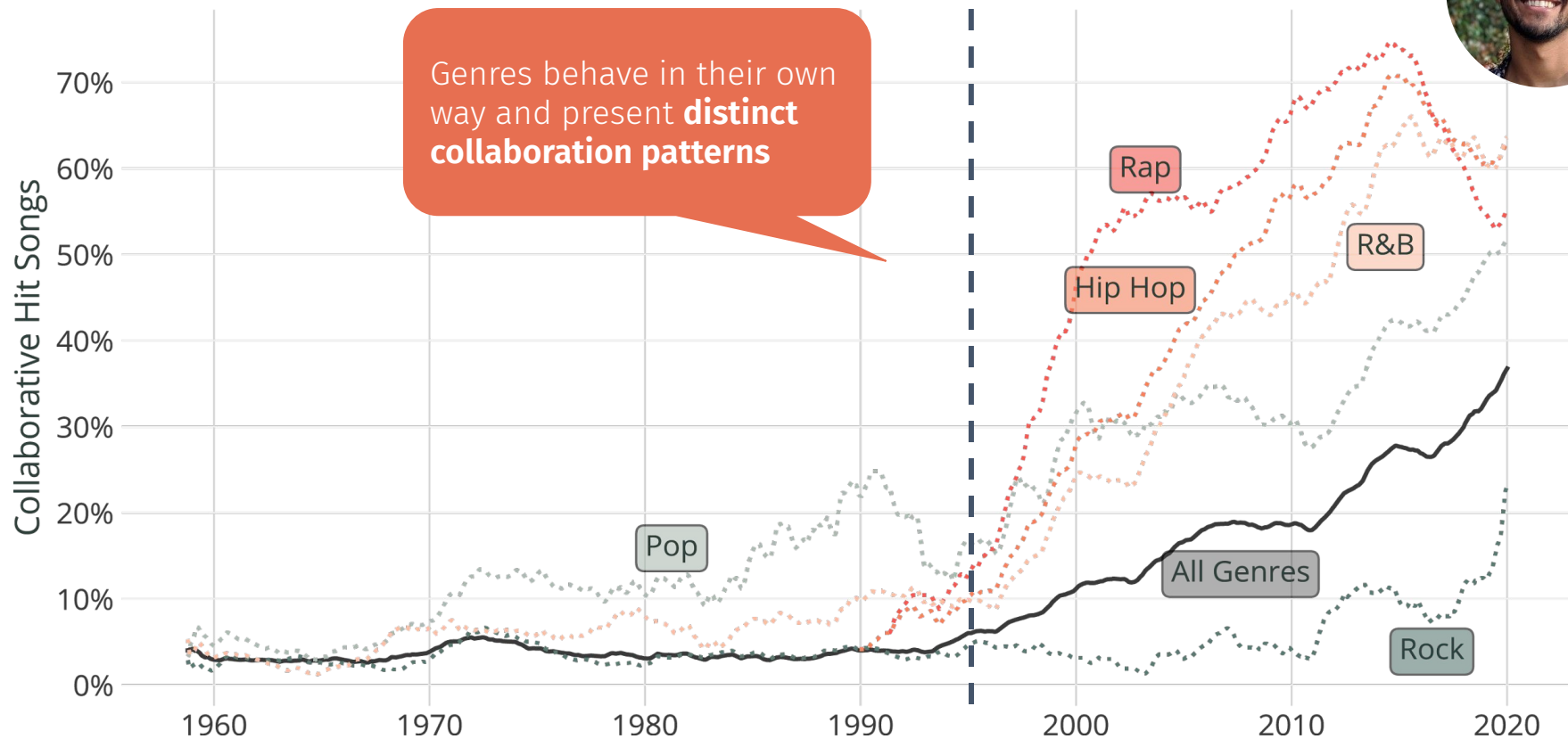
▲ Regular

● Absent



Collaboration is a powerful way to reach new audiences and increase popularity







Brazilians have a strong preference for local artists regardless of era



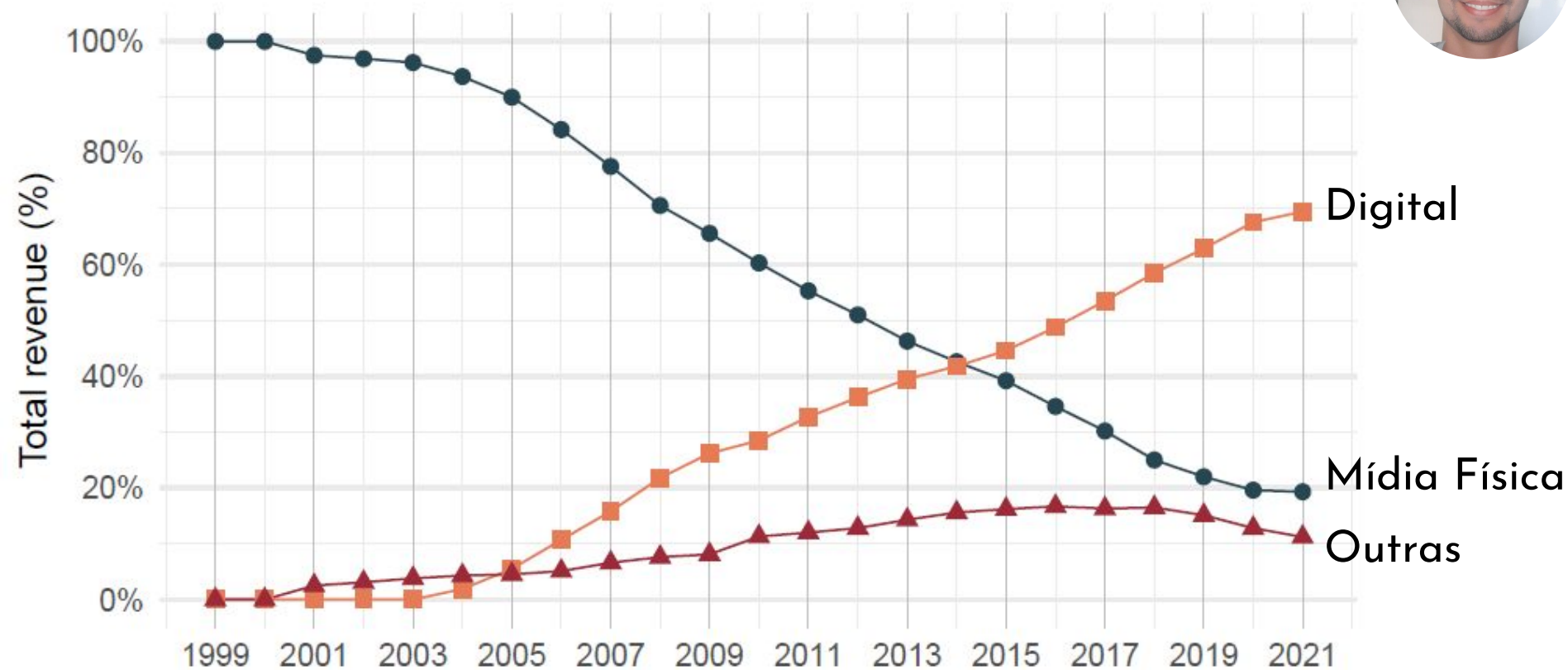
SBCM 2020

Best

Hot Streaks in the Brazilian Music Market: A Comparison Between Physical and Digital Eras

Gabriel R. G. Barbosa, Bruna C. Melo,
Gabriel P. Oliveira, Mariana O. Silva,
Danilo B. Seufitelli

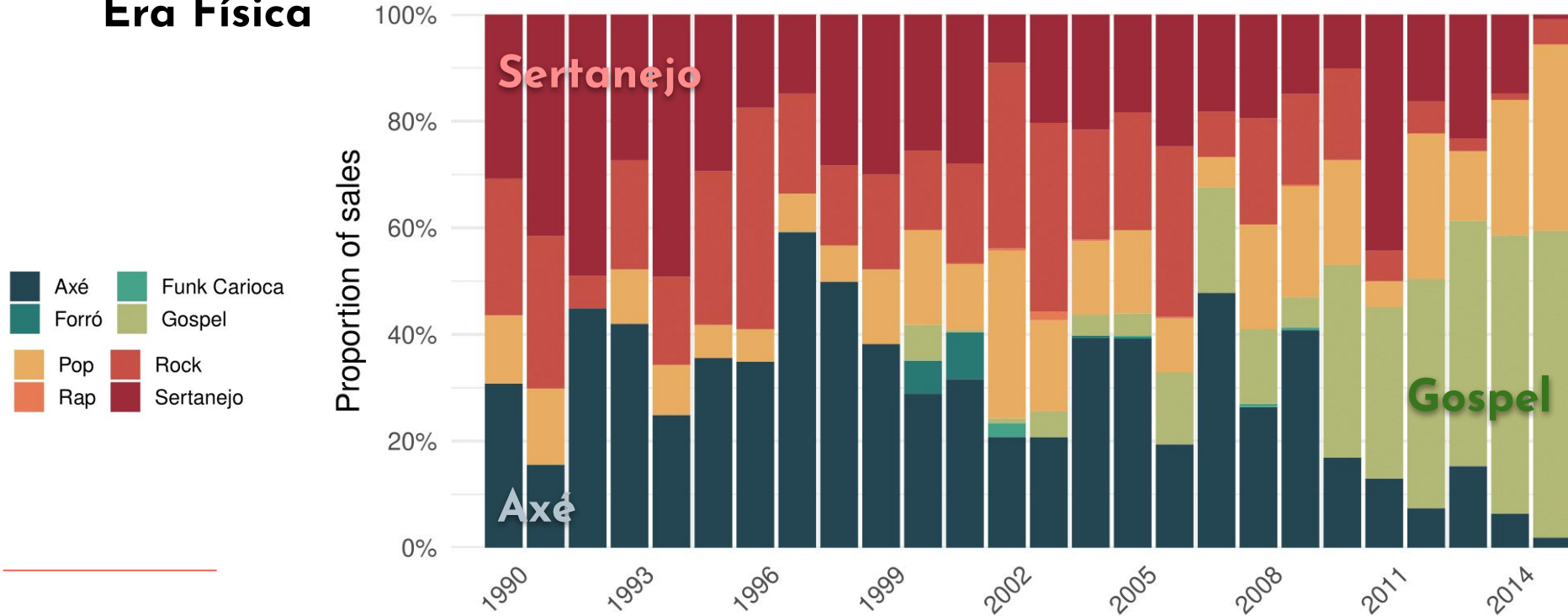
Vendas Globais





Evolução Gêneros

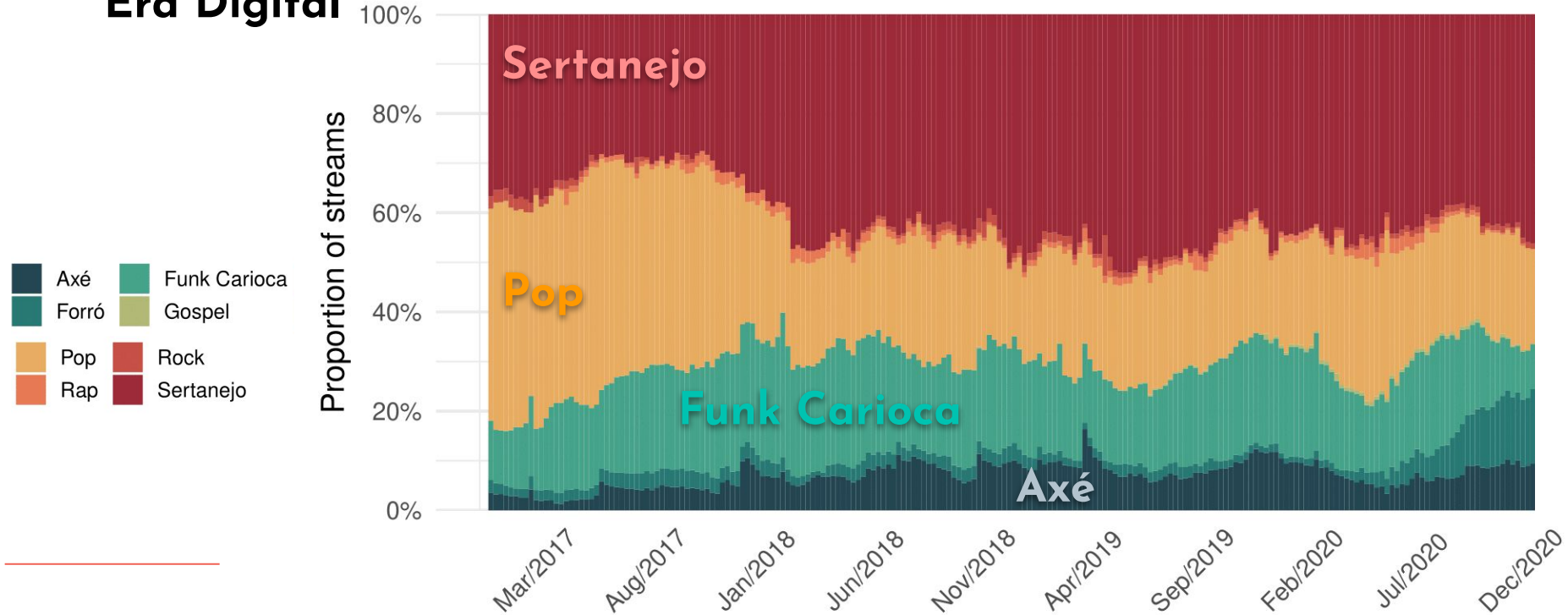
Era Física





Evolução Gêneros

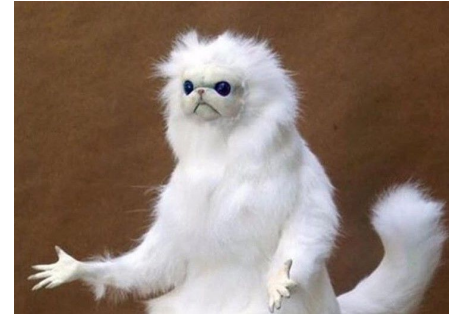
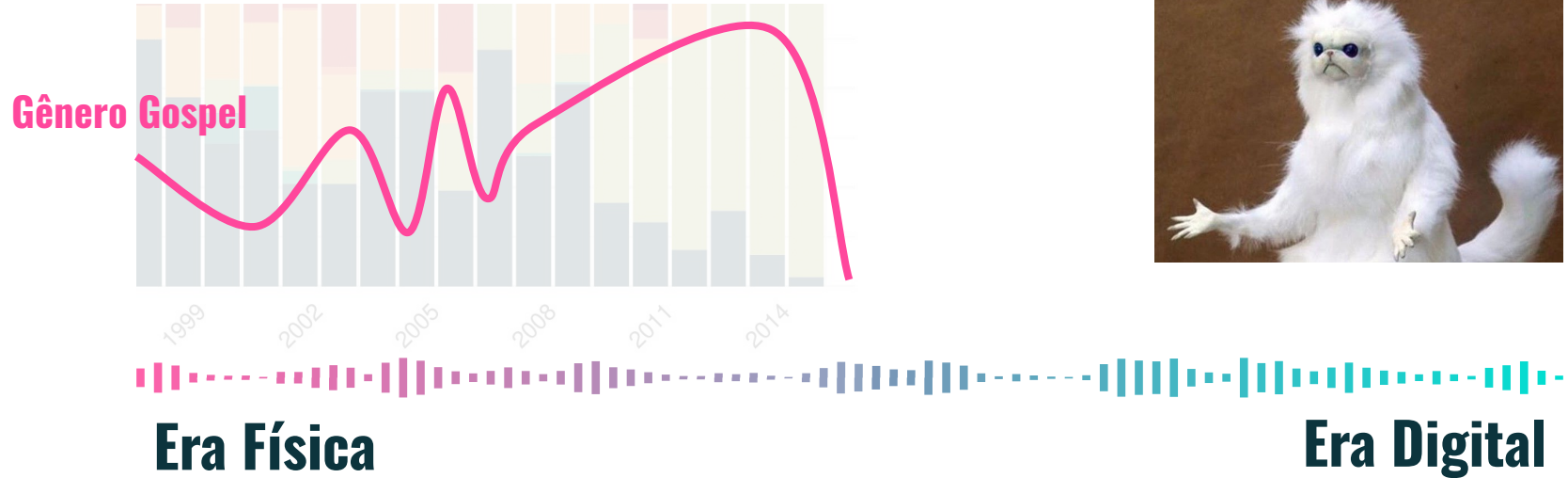
Era Digital





Evolução Gêneros

Consumo do Gênero Gospel na Era Digital



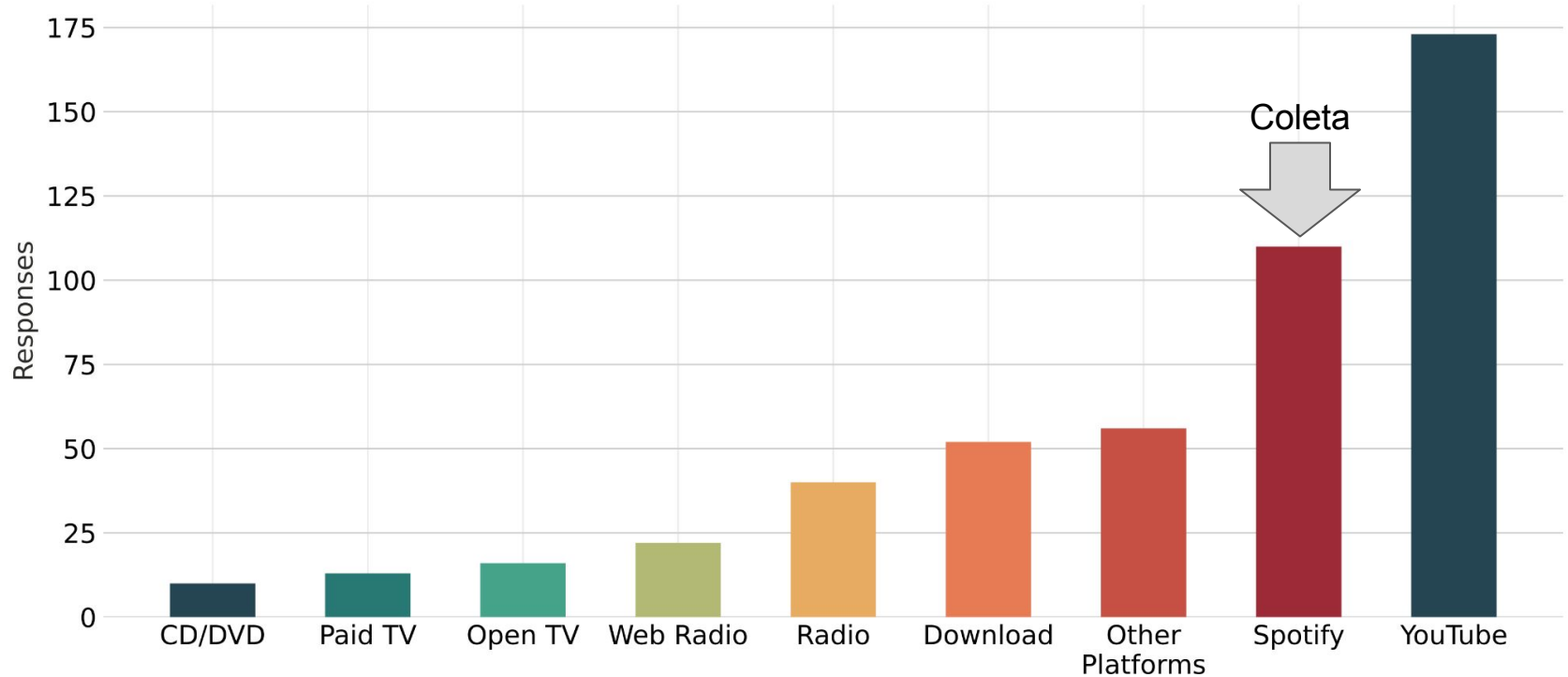
**“Algo inusitado
nos dados? ‘Bora
lá investigar”**





Evolução Gêneros

Consumo do Gênero Gospel na Era Digital



Sucesso + Gênero + Colaborações

“Exploration to Exploitation”: **SBBD 2023**

Sucesso@ Tempo, Gênero (musical), Colaboração

WTDBD @ SBBD 2020

Short @ SBBD 2022

CTDBD 2023

Datasets

DSW@SBBD 2019 MusicOSet

DSW@SBBD 2021 → **JIDM 2022** MUHSIC

DSW@**SBBD 2023** MGD+



Juliana

Sucesso + Gênero + Colaborações

Rede de Colaboração: ACM SAC 2019

WebMedia 2019

ISMIR 2020 (Best Paper Presentation)

Hot Streaks @Brazil: SBCM 2021 (Best Paper) → **Vórtex 2022**

@US: **Scientometrics 2023** (sept.23)

Hit Song Prediction: WebMedia 2022 → **JIS 2023**

Análise de Humor: CTIC @ WebMedia 2022

Descoberta de Padrões: BraSNAM 2023; **KDMiLe 2023**

Temporal Analyses: **Vórtex 2023**

Survey on Hit Song Science: **Journal 2023?** (forever!)

**“Fazer ciência é
muito mais do
que publicar
artigo científico”**



Ciência Aberta + Ensino

CoMusic - 2019 @ [Zenodo](#)

MusicOSet - 2019 @ [Zenodo](#)

Music Genre Dataset (MGD) - 2020

2855 views and 1869 downloads @ [Zenodo](#)

Music-oriented Hot Streak Information Collection (MUHSIC) - 2021

425 views and 88 downloads @ [Zenodo](#)

MGD+ 2023 @ [Zenodo](#)

Ciência de Dados usando Jupyter (com nossos Datasets)

Capítulo de Livro + Minicurso @ JAI 2021

Minicurso @ [SBBD 2021](#), Escola de Inverno PPGC/UFF

(Ciência de Dados)

Computação e Cultura



PPORTAL:

Public Domain
Portuguese-language
Literature Dataset

Mariana O. Silva	mariana.santos@dcc.ufmg.br
Clarisse Scofield	clarissescofield@dcc.ufmg.br
Mirella M. Moro	mirella@dcc.ufmg.br



Best

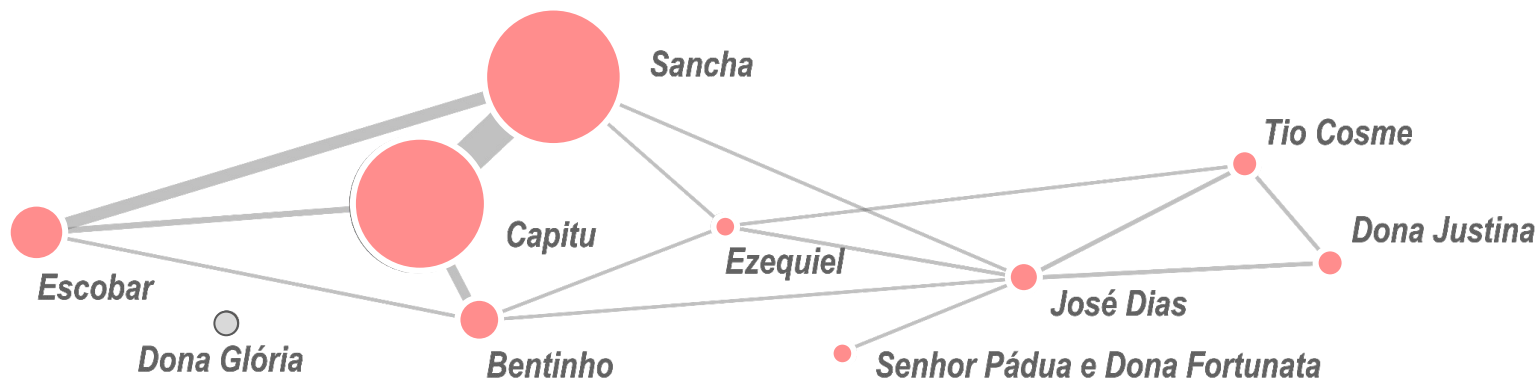
DSW
2021



CHARACTER NETWORK - DOM CASMURRO

● **Nodes** → characters

— **Edges** → interaction





Mariana, Clarisse, Luiza, Danilo, Gabriel

Ciência de Dados e Literatura

Gender Representation in Literature: Analysis of Characters' Physical Descriptions - KDMiLe 2023

Analyzing Character Networks in Portuguese-language Literary Works - BraSNAM 2023

Book Genre Classification Based on Reviews of Portuguese-Language Literature - PROPOR 2022

Cross-collection Dataset of Public Domain Portuguese-language Works - JIDM 2022

PPortal: Public Domain Portuguese-language Literature Dataset - SBBD DSW 2021 BEST

Brazilian Reading Preferences in Goodreads: Cross-state and Cross-region Analyses - iSys 2022

Exploring Brazilian Cultural Identity Through Reading Preferences - BraSNAM 2021 BEST Runner Up



Dilúvio de Dados

Uma Boia à Vista,
ou seria um Cruzeiro?

**Ao mesmo
tempo que...**

Engenharia e Ciência de Dados



Carolina Bigonha, 2012
Análise de Sentimento



Pedro O Santos, 2015
Geo + NoSQL



Levy Silva, 2018
Indexação + Deduplicação



Harley Lima, 2014
Coautoria



Luciana Maroun, 2016
Ranking



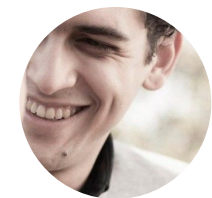
Natália Machado, 2018
Recomendação



Wladston Ferreira, 2015
Cinema



Thiago Prado, 2017
Recomendação POI



André Gonzaga, 2019
Knowledge Graphs

Pôster/Short@SBBD

CTDBD

Full e Short@SBBD

**“Fazer ciência é
muito mais do
que publicar
artigo científico”**



Banqueira de Dados

2002-2006 Revisor Externo SBBD

2008-atual Membro Comitê de Programa SBBD

Membro Comitê de Programa SBBD Componentes (NEXT)

Banca: WTDBD n vezes, CTDBD 2023, Best Paper

2010 Coordenadora Organização

2010-2012 Editora Associada JIDM

10/2013 - 10/2017 Comitê Diretivo do SBBD / CEBD

2013 Coord. Tutoriais

2014 Coord. Comitê de Programa

2015 Coord. CEBD

2017 (Proponente), 2019, 2021-2023 Dataset Showcase Workshop

2008 Coord. Sessão de Pôsteres

2009 Demos Banca Avaliadora

Banqueira de Dados

Comitê de Programa

	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23
Full		X	X		X	X	X	X	X	X	X	X	X	X	X	X	X
Pôster Short	X	X	X				X		X	X	X	X	X	X	X	X	X
Demo App	X	X	X		X				X		X	X	X				
WTDBD			X		X	X	X	X			X	X	X				X
CTDBD									X		X		X		X		X
DSW											X		X		X	X	X

Banqueira de Dados

Submissão e Autoria (*Rejeição*+Publicação)

	07	08	09	10	11	12	13	14	15	16	17	18	19	20	21	22	23
FULL		2		1		1			1+1		2	1	1		1		1+1
Pôster Short			1	1	2+2	4	2+1	1	3	4	1+4	1			2	1	
Demo App				1	1		2		4		1						
WTDBD		1										1	1	1			
WTAG																1	
CTDBD	-	-	-	-	-	-	-	-	MM	-		-	d+M		MM		M
DSW													1		2		1

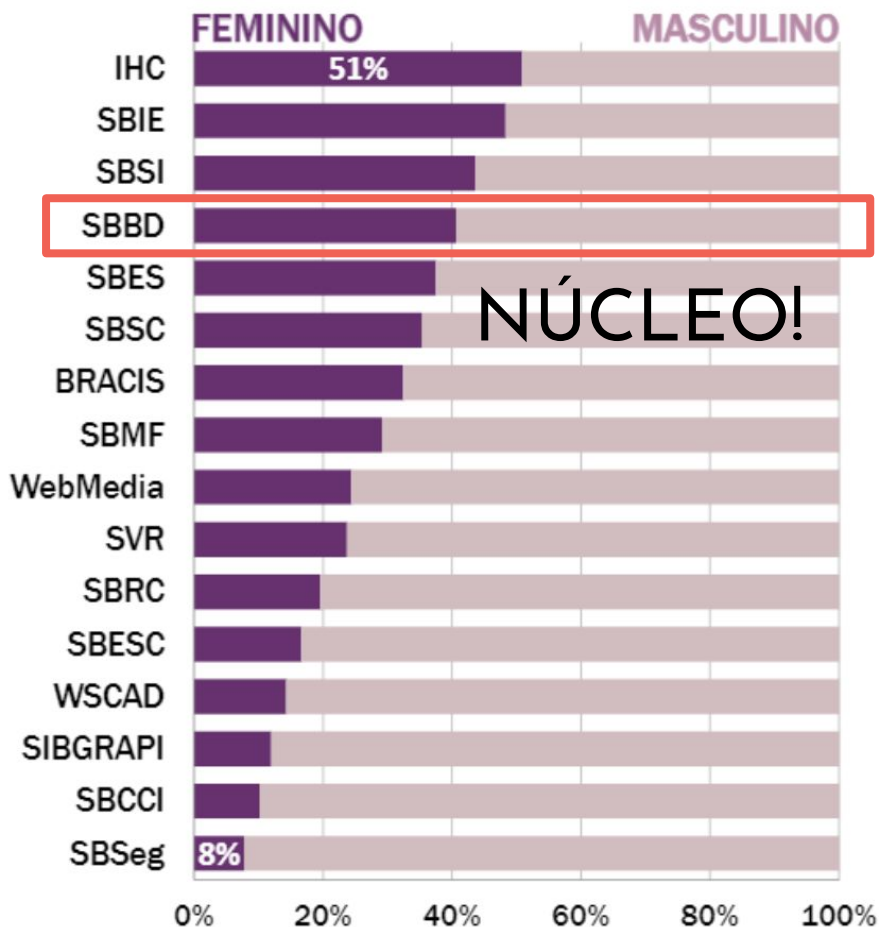
Banqueira de Dados

BEST

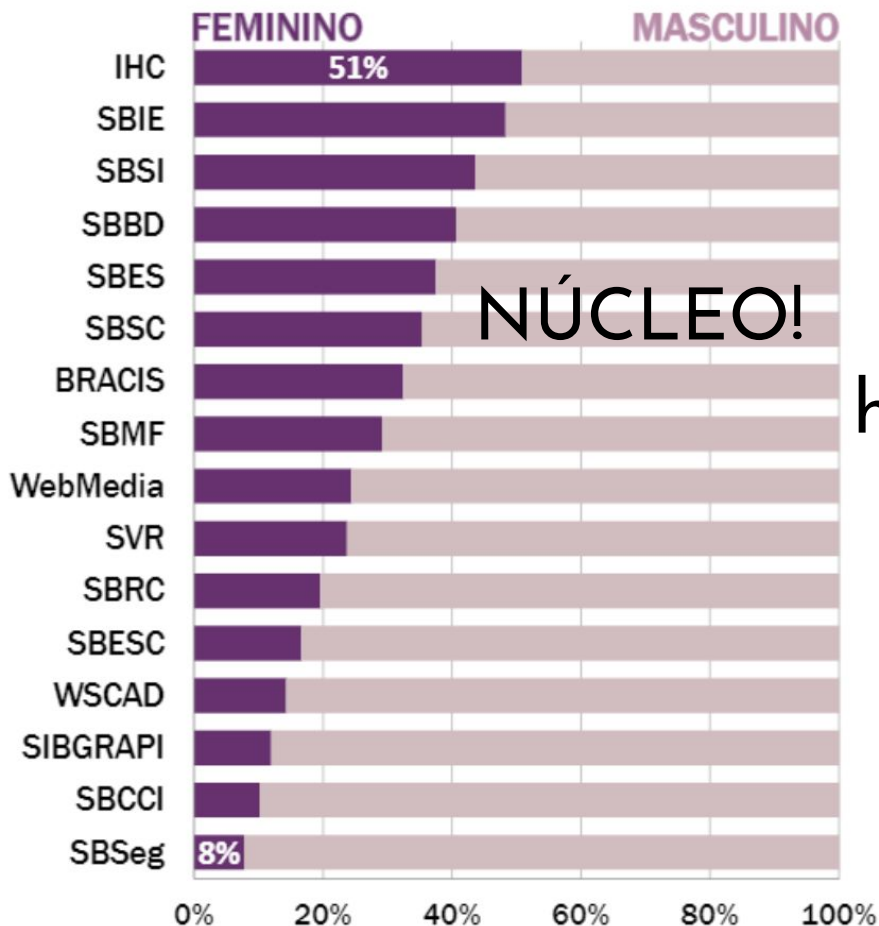
- > 2023 MH DSW: MGD+
- > 2021 DSW: PPORTAL
- > 2017 Short: GitHub
- > 2017 Runner Up: STACY
- > 2013 Short@JIDM
- > 2012 Short@JIDM
- > 2011 Menção Honrosa Demo
- > 2010 Short@JIDM
- > Best Reviewer



Diversidade & Computação



Diversidade & Computação



<https://linktr.ee/lackofdiversity>

Best @ EduComp 2022

Women in dataBases





2008
foto
Campinas



2009 Fortaleza



**2010
BH**



2011 Florianópolis





2012
São
Paulo





2014
Curitiba



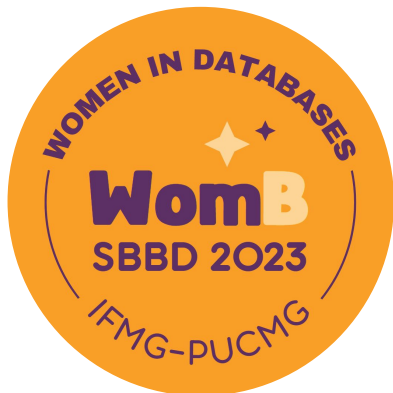
2015 Petrópolis

**2016
Salvador**



2022 Buzios





#gratidão



Dilúvio de Dados

Uma Boia à Vista, ou seria um Cruzeiro?

MIRELLA M. MORO

UFMG

bit.ly/profamirella



**SBBD
2023**

Créditos

Imagens: PixaBay

Fotos: acervo pessoal

Template: Macari Company Profile by Slidesgo