# Aprendizagem de máquina aplicada a indexação automática de documentos digitais

**III Seminário do Portal de Periódicos - CAPES**

**Dalton Martins**

Faculdade de Ciência da Informação

Universidade de Brasília

daltonmartins@unb.br

FACULDADE DE CIÊNCIA DA INFORMAÇÃO

UnB, sua linda
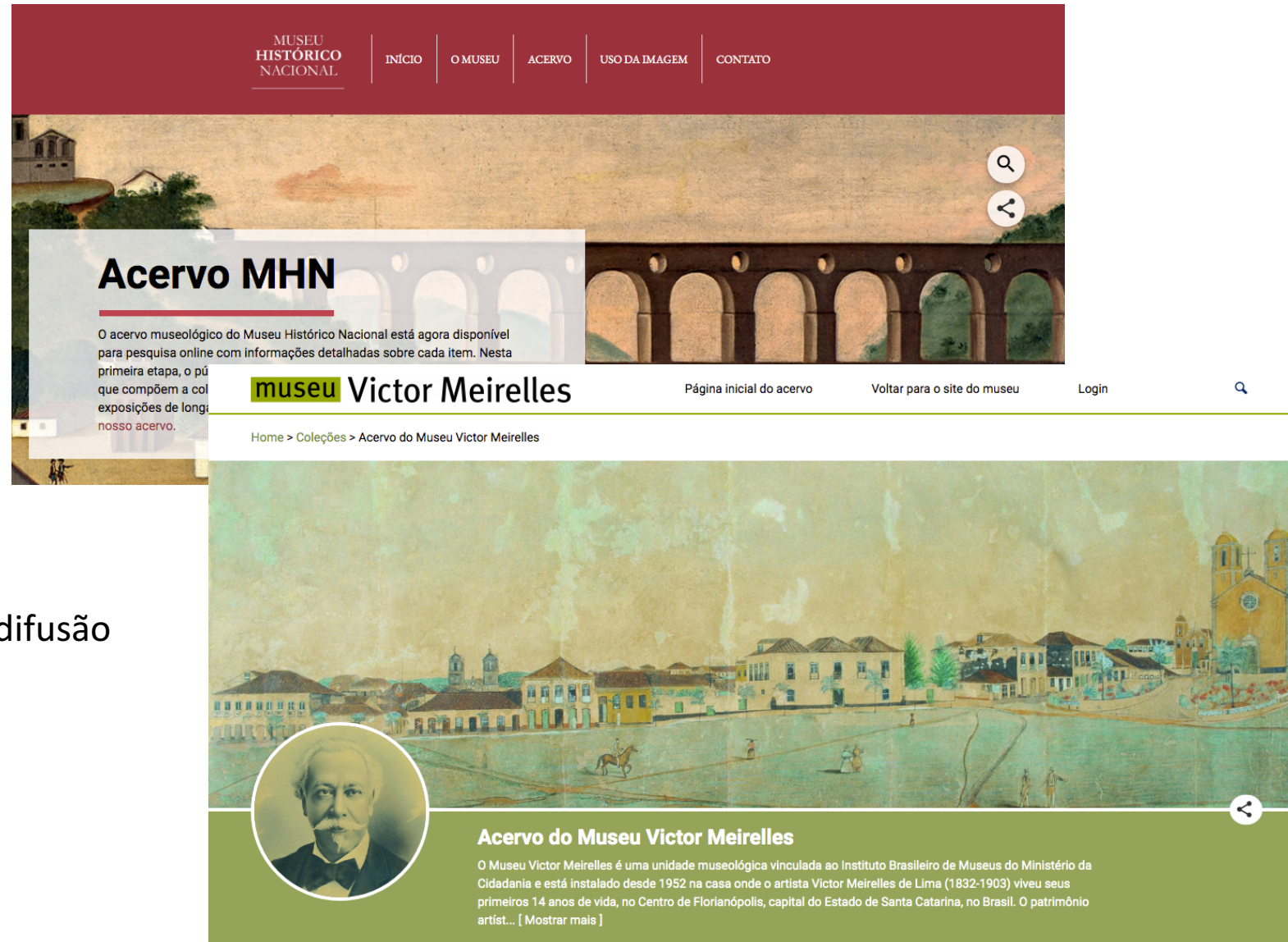
meu orgulho é você

# Quem somos nós: institucionalidade



Universidade de Brasília
Faculdade de Ciência da Informação



**Laboratório de Inteligência de Redes**
Biblioteca Central

# Quem somos nós: principais projetos



- **Pesquisa e desenvolvimento do Tainacan:**
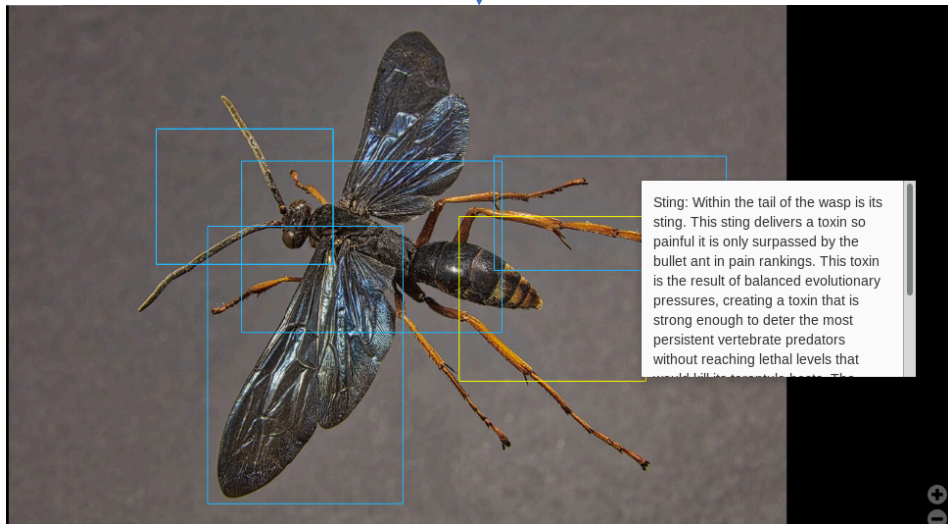repositório digital para organização, gestão e difusão de acervos digitais em rede.

# Quem somos nós: principais projetos

- Reconhecimento e reconciliação semântica de entidades em texto;
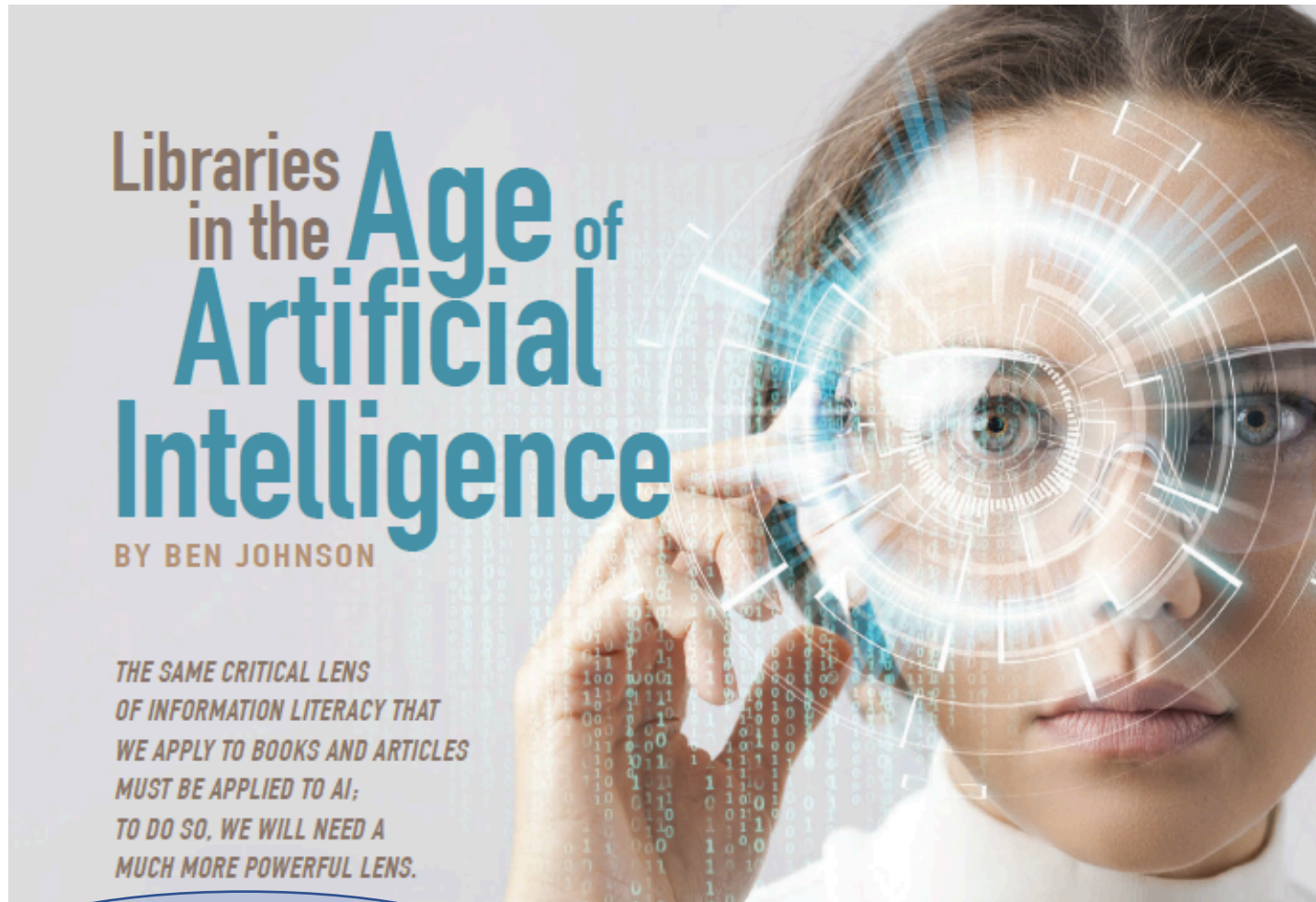- Indexação automática e semi-automáica em imagens.

Pela primeira vez documentada no século XIII, Berlim foi sucessivamente a capital do Reino da Prússia (1701), do Império Alemão (1871-1918), da República de Weimar (1919-1932) e do Terceiro Reich (1933-1945). Depois da Segunda Guerra Mundial, a cidade foi dividida. Berlim Oriental se tornou a capital da República Democrática Alemã (RDA), enquanto Berlim Ocidental continuou sendo parte da República Federal da Alemanha (RFA).18 Com a reunificação alemã em 1990, a cidade passou a ser capital de toda a Alemanha.

Sting: Within the tail of the wasp is its sting. This sting delivers a toxin so painful it is only surpassed by the bullet ant in pain rankings. This toxin is the result of balanced evolutionary pressures, creating a toxin that is strong enough to deter the most persistent vertebrate predators without reaching lethal levels that

Pela primeira vez documentada no século XIII, Berlim foi sucessivamente a capital do Reino da Prússia (1701), do Império Alemão (1871-1918), da República de Weimar (1919-1932) e do Terceiro Reich (1933-1945). Depois da Segunda Guerra Mundial, a cidade foi dividida. Berlim Oriental se tornou a capital da República Democrática Alemã (RDA), enquanto Berlim Ocidental continuou sendo parte da República Federal da Alemanha (RFA).18 Com a reunificação alemã em 1990, a cidade passou a ser capital de toda a Alemanha.
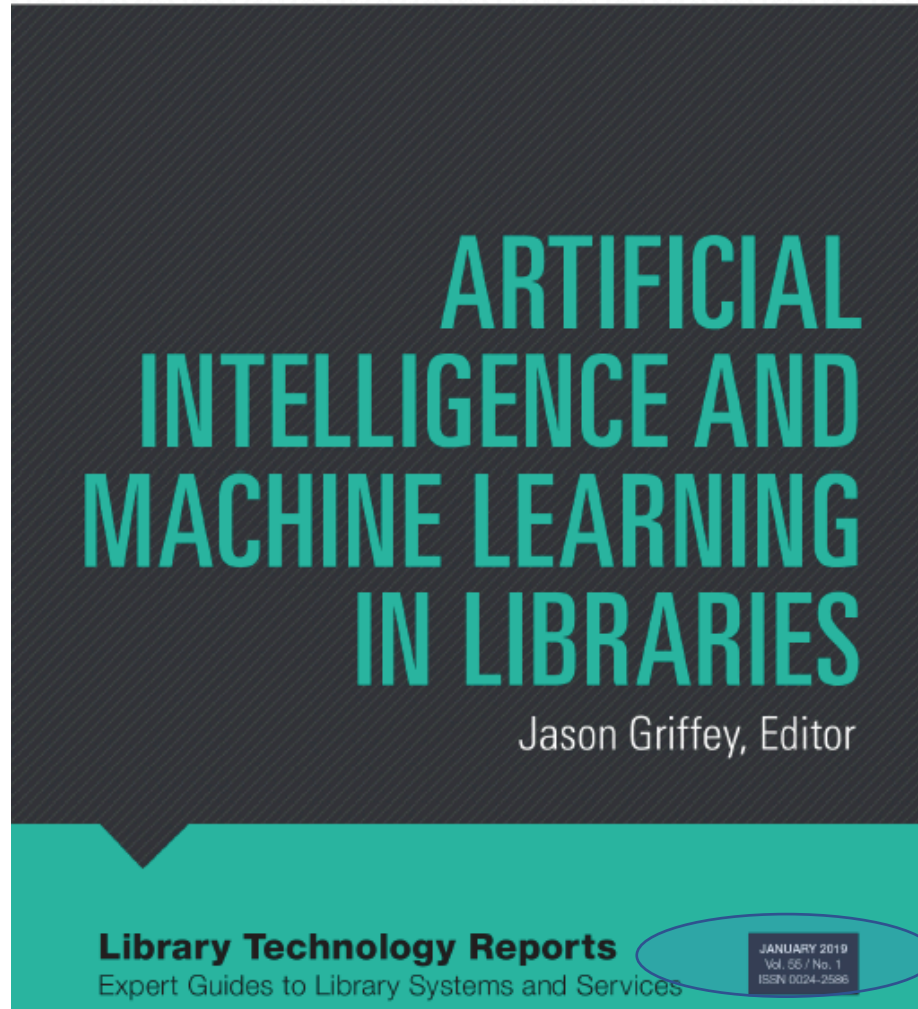
# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente



Revista InfoToday
2018

# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente

ALAAmericanLibraryAssociation

ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING IN LIBRARIES

Jason Griffey, Editor

**Library Technology Reports**
Expert Guides to Library Systems and Services

JANUARY 2019
Vol. 55 / No. 1
ISSN 0024-2586

American Library Association
2019
Caderno especial de tecnologia para as bibliotecas.

# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente



American Library Magazine
2019
Editorial especial

# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente



American Library Magazine
2019
Análise de tendências

# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente



*Public Libraries Leading the Way*

**The Democratization of Artificial Intelligence: One Library's Approach**

Thomas Finley

Chances are that before you read this article, you probably checked your email, used a mapping app to find your way, or typed a search term online. Without your even perceiving it, artificial intelligence (AI) has already helped you to accomplish something today. Email spam filters use variants of AI to help cut down on harmful or useless emails in your inbox.[1] With AI doing the fact-crunching, mapping apps quickly preview the best route based on a myriad of factors. Search engine companies like Google have been using AI to suggest or produce results faster for longer than anyone outside of the company really knew until recently.[2] According to a recent study by Northeastern University and Gallup, 85% of Americans are already using AI products.[3] The true revelation behind these recent technological developments may not be the fact that AI is already embedded into the fabric of our modern lives. The real surprise might just be the sudden ubiquitous availability (and approachability) of AI tools for all. As Google's former Chief Scientist of AI and Machine Learning, Fei-Fei Li, said in 2017, "The next step for AI must be democratization, lowering the barriers of entry, and making it available to the largest possible community of developers, users and enterprises."[4] This sounds a lot like most public libraries' mission statements. As with other important workforce development efforts, libraries are uniquely placed to participate in this new revolution as key platforms for the discovery and dissemination of emerging tech knowledge. At the Frisco Public Library (https://www.friscolibrary.com), we saw this AI trend surfacing, we see AI as a critical future job skill, and we investigated ways to introduce our patrons into this space. As such, the Frisco Public Library has leveraged readily available technology in a cost-effective way that has engaged community interest. Our efforts are also replicable and scalable in terms of multi-nodal experiences both at home and in classroom-based learning.

**SOME BASIC DEFINITIONS**

Let's take a few steps back to give some broad definitions and boundaries to the scope of AI. According to the Oxford English Dictionary, artificial intelligence is "the capacity of computers or other machines to exhibit or simulate intelligent behavior."[5] In the literature, you will find a further distinction between General AI, Narrow AI, and something called Machine Learning.[6]

General AI is something that begins to look like science fiction: an artificial intelligence that learns how to learn, then is able to generalize what it has learned and apply that knowledge to a different case. In advanced examples of General AI, scientists are thinking of not putting a specific problem in front of a General AI program to solve, rather, they are giving it an entire dataset so the program *itself* can choose what problems it should work on. Removing the limited point of view of whoever programs the program.[7]

Narrow AI is easier to understand because it is what we interact with the most in our day-to-day lives. It is what powers those little speed ups that help us do things faster every day: search

**Thomas Finley** (tfinley@friscotexas.gov) is Adult Services Manager, Frisco Public Library.

INFORMATION TECHNOLOGY AND LIBRARIES | MARCH 2019                    8

Information Technology and Libraries
2019
Revista científica

# Aprendizagem de máquina e inteligência artificial aplicadas às bibliotecas: um tema emergente



Trabalhos apresentados a IFLA "inteligência artificial"

Trabalhos apresentados nos congressos da IFLA
Descritor: "inteligência artificial"

Há um crescente e importante interesse nos temas aprendizagem de máquina e inteligência artificial no campo das bibliotecas nos últimos anos.

O que isso significa em termos de oportunidades de novos serviços, melhoria na qualidade das instituições e necessidade de qualificação profissional?

# O que os bibliotecários estão pensando sobre isso?

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA

Pesquisa realizada durante os meses de maio e junho de 2017 via listas de emails de profissionais bibliotecários nos EUA.

Publicada na revista InfoToday em fevereiro de 2018

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA



In your opinion, what departments in the library are most likely to be affected by Artificial Intelligence

Virtual Services/Discovery — 247 — 77.92%

Reference — 222 — 70.03%

Cataloging — 164 — 51.74%

Collection Development — 107 — 33.75%

Acquisitions/Technical Services — 105 — 33.12%

Instruction — 98 — 30.91%

Access Services — 83 — 26.18%

Administration — 33 — 10.41%

Chart 2: question two (n=317)

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA



**Within 2 (2019), 10 (2027), and 30 (2047) years what is the probability of supercomputers like Watson being used in the library?**

Chart 3: questions three (n=319), four (n=318), and five (n=317)

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA



How positive or negative will the overall impact of supercomputers like Watson be on the profession of librarianship?

NEGATIVE. 27.36%
NO EFFECT. 15.96%
POSITIVE. 56.68%

Charts 4, 5, and 6: questions six (n=315), seven (n=315), and eight (n=307)

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA



Chart 7: question nine (n=309)

# A percepção dos bibliotecários sobre a aprendizagem de máquina e IA



Chart 8: question 10 (n=310)

# Como esse cenário impacta as bibliotecas?

# Dados de qualidade



A grande riqueza dos dados das bibliotecas está em suas coleções, na catalogação, na indexação, na relação com o uso da informação!

São insumos fundamentais para a geração de novos serviços e produtos informacionais!

Novos serviços que usem esses dados são um próximo estágio fundamental ao desenvolvimento das bibliotecas.

# Tipos de serviços
## (mais rapidamente impactados pela IA)

- Referência
- Catalogação e indexação automática ou semi-automática
- Busca e recuperação da informação
- Aquisição
- Processamento de linguagem natural, imagem, áudio e vídeo
- Reconhecimento de padrões e redes semânticas de relacionamento entre documentos.

# Potencial de aplicação da IA em Bibliotecas

# Potencial de aplicação da IA em Bibliotecas

| Library roles in AI | Competencies that need to be extended | Alternative providers of service/ function |
|---|---|---|
| Procuring content for AI to work from (including both licensing and through open access) | Procurement and licensing of e-content | Publishers and other new intermediaries |
| Providing content | Digitisation, metadata provision | Publishers and other new intermediaries |
| Data quality control | Collection management | |
| Procuring AI tools | Procurement and licensing of software and services | IT departments, academic departments |
| Data curation (e.g. of derived data) | Collection management, digital preservation | Publishers and other new intermediaries |
| Designing data infrastructure to enable AI | Design of information discovery infrastructure | IT departments |
| Explaining how to navigate the new information environment | Understanding of the scholarly publishing landscape, including data creation processes | |
| Teaching critical data literacy: understanding how to evaluate AI tools and their results, and also protect one's own privacy | Information literacy | IT departments |
| Designing AI tools | N/A – outside normal library professional work | Academic departments, Publishers |
| Data analysis and writing algorithms | N/A | IT departments, academic departments |

Table 1 Potential library roles in AI

# Impactos em potencial para as bibliotecas

Others can only be guessed at. Yet bringing our data together with the literature, it emerges that AI should be seen as defining a nexus of change (Pinfield et al., 2017) that has "wide and deep" ramifications (JISC reference) in terms of:

1. What a library is, what a collection is and how to search for material. The library may increasingly be seen as data, accessed through AI, the scope of the collection as framed by the AI;
2. How established services are delivered, for example by chatbots and other intelligent agents;
3. What users expect of libraries: through expectations learned in other areas;
4. What information literacy is: the ability to navigate a new space of AI tools and data, and data literacies, including critical awareness of how to protect one's own privacy;
5. Who users are: some users will be AI tools; human access to content will be remediated through content being summarised and partially analysed for them by machines;
6. What libraries know about users and so how the library is managed: because of management decisions based on use data, combined with other learning and research analytics;
7. How the library works with other internal and external partners and competitors, especially IT services and new third-party commercial services;
8. How library services are evaluated: again through wider and deeper data;
9. What skills librarians need: be that for licensing, evaluation of data analysis and visualisation tools or using such tools themselves;
10. Whether the library community can operate effectively at different levels beyond the institution: in order to design and deliver services which will serve international communities of scholars and students;
11. Indeed, whether we need librarians (because of chatbots, automated metadata creation tools etc) or libraries (because of alternative intermediaries) at all, at least as currently conceived.

# Exemplos e aplicações

drop image
or
click here

akiwi finds keywords for your images.

http://www.akiwi.eu/

# Alguns serviços
# (exemplos)



https://hamlet.andromedayelton.com/

# Alguns serviços
# (exemplos)

## API CLOUD VIDEO INTELLIGENCE

Pesquise e descubra seu conteúdo de mídia com a API Cloud Video Intelligence

FAÇA UMA AVALIAÇÃO GRATUITA

## Análise precisa de vídeo

Com a API Google Cloud Video Intelligence, a pesquisa e a descoberta de vídeos ficaram ainda melhores ao coletar metadados com uma API REST muito fácil de usar. Agora é possível pesquisar cada momento de cada arquivo de vídeo no seu catálogo. Com essa API, você pode fazer anotações em vídeos armazenados no Google Cloud Storage

https://cloud.google.com/video-intelligence/

# Alguns serviços
# (exemplos)

## Análise avançada de imagens

O Cloud Vision oferece modelos pré-treinados por meio de uma API ou a capacidade de criar modelos personalizados usando o AutoML Vision. Assim, você tem a flexibilidade que precisar, dependendo do seu caso de uso.

A **API do Cloud Vision** encapsula modelos avançados de machine learning em uma API REST fácil de usar, o que permite aos desenvolvedores entender o conteúdo de imagens. Essa API rapidamente classifica as imagens em milhares de categorias, (por exemplo: "veleiro"), detecta objetos e rostos individuais e extrai palavras impressas contidas nas imagens. Crie metadados no seu catálogo de imagens, modere conteúdo ofensivo ou ative novos cenários de marketing usando a análise de sentimento das imagens.

https://cloud.google.com/vision/

# Uma agenda de P&D para o desenvolvimento da IA nas bibliotecas

- Formação:
  - Incorporar disciplinas e projetos de pesquisa na formação de graduação dos bibliotecários;
  - Capacitar profissionais atuando na área em cursos de média e longa duração, como especialização;
- Desenvolvimento:
  - Diagnosticar conjuntos de dados que podem ser abertos para novos serviços pelas bibliotecas;
  - Implementar descrição semântica dos dados para melhorar qualidade descritiva e facilitar consumo por algoritmos computacionais;
  - Selecionar e priorizar serviços de maior interesse para a geração de novos produtos informacionais automatizando o acesso aos dados;
- Pesquisa:
  - Incorporar algoritmos de aprendizagem de máquina na análise e tratamento de dados oriundos das bibliotecas;
  - Propor novos modelos conceituais para automatizar serviços de busca, recuperação, indexação, catalogação e descoberta de novos conhecimentos nos conjuntos de dados.
  - Realizar projetos experimentais com as instituições de interesse.

# Obrigado!

## daltonmartins@unb.br