

Uma estratégia para coleta, integração e tratamento de dados científicos no contexto do BrCris

A strategy for collection, integration, and processing of scientific data in BrCris context

Washington Segundo

IBICT, Brasil.

E-mail: washingtonsegundo@ibict.br

Thiago Magela Dias

CEFET-MG, Brasil.

E-mail: thiagomagela@cefetmg.br

Tales Moreira

CEFET-MG, Brasil.

E-mail: tales.info@gmail.com

Adilson Luiz Pinto

Universidade Federal de Santa Catarina, Brasil.

E-mail: adilson.pinto@ufsc.br

Vivian Silva

IBICT, Brasil.

E-mail: vivian.ss@gmail.com

Josir Gomes

IBICT, Brasil.

E-mail: josir@irdx.com.br

Luc Quoniam

UFMS, Brasil.

E-mail: quoniam.luc@gmail.com

Como citar: Segundo, W.; Dias, T. M.; Moreira, T.; Pinto, A. L.; Silva, V.; Gomes, J.; Quoniam, L.; Matas, L.; Dias, A.; & Schneider, J. (2022). Uma estratégia para coleta, integração e tratamento de dados científicos no contexto do BrCris. En T. M. R. Dias (Ed.), *Informação, Dados e Tecnologia. Advanced Notes in Information Science, volume 2* (pp. 215-222). Tallinn, Estonia: COLNes Publishing. DOI: 10.47909/anis.978-9916-9760-3-6.117.

Copyright: © 2022, The author(s). This is an open access work distributed under the terms of the CC BY-NC 4.0 license which permits copying and redistributing the material in any medium or format, adapting, transforming and building upon the material as long as the license terms are followed.

Lautaro Matas

La Referencia, Brasil.

E-mail: lmatas@gmail.com

Ary Dias

IBICT, Brasil.

E-mail: arygabrieldias@gmail.com

Juliana Schneider

IBICT, Brasil.

E-mail: jschneider.js@gmail.com

RESUMO

Nos últimos anos várias iniciativas que visavam a criação de sistemas que gerenciam a produção acadêmica de uma instituição, país ou área de conhecimento têm recebido atenção de diversas áreas. Tais sistemas são conhecidos pela sigla CRIS (*Current Research Information Systems*) e têm como objetivo agregar informações de bases de dados diversas com intuito de fornecer relatórios e dados consolidados para que pesquisadores possam analisar. Logo, este trabalho apresenta o processo de desenvolvimento da Plataforma BrCris com o objetivo de fornecer ferramentas tecnológicas com o intuito de munir a comunidade acadêmica brasileira com dados consolidados da produção científica nacional. Tal iniciativa se apresenta como importante mecanismo de agregação de dados, fornecendo visualizações e análises de importantes conjuntos de dados científicos em especial sobre a produção científica brasileira. Logo, a disponibilização da Plataforma BrCris para a comunidade científica irá proporcionar diversos estudos bibliométricos que a priori seriam de extrema complexidade na sua concepção.

Palavras-chave: Produção Científica; BrCris; Plataforma Lattes.

ABSTRACT

Several initiatives to create systems that manage the academic production of an institution, country, or area of knowledge have received attention from different fields. Such systems are known by the acronym CRIS (Current Research Information Systems) and aim to aggregate information from other databases to provide reports and consolidated data for researchers to analyze. Therefore, this work presents the development process of the BrCris Platform to provide technological tools to give the Brazilian academic community consolidated data from the national scientific production. This initiative presents itself as an essential data aggregation mechanism, providing views and analysis of critical scientific datasets, especially on Brazilian scientific output. Therefore, making the BrCris Platform available to the scientific community will give several bibliometric studies that a priori would be highly complex in their conception.

Keywords: Scientific Production; BrCris; Lattes Platform.

1. INTRODUÇÃO

A PRODUÇÃO científica brasileira tem crescido expressivamente e, em perspectiva às especificidades de campos disciplinares distintos, heterogênea quanto à tipificação de sua produção tanto em termos quantitativos como qualitativos. E o resultado desta produção se materializa em forma de artigos em periódicos, teses e dissertações, além de produtos diversos como: softwares, patentes, obras e instalações artísticas, entrevistas e projetos cinematográficos.

Para o campo da Ciência da Informação, e em especial da Cientometria, quantificar essa produção é uma tarefa árdua pois a disponibilização de bases de dados abertas muitas vezes é restrita ou simplesmente inexistente em determinados contextos. Bases de dados proprietárias como a Scopus, Web of Science, Google Acadêmico e Microsoft Research Data permitem o acesso mas este é sempre limitado ao número de registros que podem ser obtidos, contemplam poucos repositórios e periódicos nacionais e ainda existe o grave problema da opacidade dos algoritmos utilizados por estas plataformas que determinam o que é ou não relevante.

A partir desse cenário, começaram a surgir iniciativas que visavam a criação de sistemas que gerenciam a produção acadêmica de uma instituição, país ou área de conhecimento. Tais sistemas são conhecidos pela sigla CRIS (Current Research Information Systems) e têm como objetivo agregar informações de bases de dados diversas com intuito de fornecer relatórios e dados consolidados para que pesquisadores da área possam analisar como se dá a produção em seus países ou áreas de conhecimento.

CRIS define um sistema de informação sobre todo o ecossistema do processo científico. São organizadas em um só lugar todas as informações do ciclo da pesquisa Científica, desde o Fomento, passando pelos Projetos, Pesquisadores, Instituições de Pesquisa e Laboratórios, até os outputs de uma pesquisa científica, tais como artigos científicos, teses, dissertações, livros, capítulos de livro, patentes e conjuntos de dados científicos (Sivertsen, 2019).

Neste contexto, a idealização do Projeto do Sistema BrCris, que é o CRIS no contexto da Ciência Brasileira, data de 2014, quando inspirado no modelo proposto por Portugal de um CRIS nacional (o PTCRIS - <https://ptcris.pt>), o Instituto Brasileiro de Informação em Ciência e Tecnologia (IBICT) iniciou uma sequência de estudos e parcerias interinstitucionais para a execução do Projeto. Em 2020, houve a implementação formal de um Projeto de Pesquisa para a construção do BrCris. O intuito é fornecer ferramentas tecnológicas visando munir a

comunidade acadêmica brasileira com dados consolidados da produção científica nacional. Tomando como base outros projetos CRIS e padrões internacionais disponibilizados pelo OpenAIRE e COAR.

Logo, o BrCris tem por objetivo estabelecer um modelo único de organização da informação científica de todo o ecossistema da pesquisa brasileira. Entre os agentes deste ecossistema estão os pesquisadores, os projetos, infraestruturas, laboratórios e instituições de pesquisa, os financiadores, além dos resultados da pesquisa expressos principalmente por publicações científicas, teses, dissertações, conjuntos de dados científicos, software e patentes.

Diante disso, com a integração dos dados em um repositório de dados padronizado, o uso do dashboard em forma de visualização elucida alguns benefícios, como a redução de complexidade de dados, auxilia na percepção das propriedades existentes, ajuda na detecção de erros aparentes, consegue englobar a representação em pouco conteúdo, amplia a percepção cognitiva, dentre outros.

Outro aspecto que vale ser salientado na função do dashboard é a geração de índice e indicadores visando atribuir a mensuração de fenômenos, sejam de natureza social, econômica ou científico/tecnológica. No contexto deste trabalho, valoriza-se o foco em sua representação e facilidade de identificação dos cenários.

Tendo em vista o que foi exposto, este trabalho tem como objetivo apresentar o processo de desenvolvimento do projeto BrCris, que visa coletar, integrar e disponibilizar informações relativas ao universo de pesquisa científica no Brasil, catalogando e traçando relacionamentos entre pesquisadores, as organizações às quais pertencem, os projetos dos quais participam e como são financiados, e todos os produtos por ele gerados, como publicações, patentes e software.

2. PROCEDIMENTOS METODOLÓGICOS

O BrCris concentra um amplo ecossistema de dados, de diversas fontes, como por exemplo, dados curriculares de indivíduos, sobre organizações, programas de pós-graduação, publicações, orientações, revistas científicas, dentre outros, sendo necessário todo um esforço para tratamento dos dados de interesse. Neste contexto, tendo em vista as diversas fontes de dados que irão compor o BrCris se faz necessário a transformação dos dados em formato padronizado, sendo necessário a transformação baseada em um modelo que será importado para a plataforma LA Referencia.

Em se tratando do modelo de dados do BrCris, iniciou-se pela adoção de nove entidades de dados, seguindo padrões amplamente utilizados na comunidade científica internacional. São elas:

- Project: projetos de pesquisa executados, ou em execução;
- Service: revistas científicas, repositórios digitais, bibliotecas digitais e outras fontes de informação científica;
- Program: programas de pós-graduação brasileiros;
- Course: cursos nacionais, ou internacionais de pós-graduação stricto ou lato sensu;
- OrgUnit: instituições, faculdades, departamentos de pesquisa;
- Person: pesquisadores, assistentes de pesquisa e pessoas de apoio técnico à pesquisa;
- Patent: patentes como resultado da pesquisa;
- Dataset: conjuntos de dados de pesquisa coletados por pesquisadores e demais agentes no âmbito de um projeto ou pesquisa científica;
- Publication: artigos científicos, teses, dissertações, livros, capítulos de livro e relatórios científicos.

O modelo de dados é definido por um conjunto de entidades e relações, que por sua vez possuem identificadores e atributos pré-definidos. A utilização de um descritivo visa facilitar a identificação de atributos de cada entidade e suas relações, possibilitando que com o auxílio de uma rotina desenvolvida especificamente para esta funcionalidade, o modelo possa incorporar todas as mudanças realizadas diretamente no modelo. Esta estratégia visa facilitar de forma significativa a incorporação de novos atributos e relações, sem a necessidade de alterações diretamente no modelo de dados.

O descritivo de uma Entidade apresenta inicialmente o atributo definido no Modelo de Dados bem como sua respectiva descrição. Logo, para cada conjunto de dados que é fonte de informações para a Entidade, são descritos os atributos dos conjuntos de dados relacionados aos do modelo.

Para o tratamento dos dados foi desenvolvida uma biblioteca na linguagem de programação Python, contendo uma estrutura de dados preparada para facilitar o processamento de dados originários de todas as fontes para o formato exigido pela plataforma LA Referencia. Logo, a biblioteca desenvolvida é responsável por toda a transformação e exportação dos dados, utilizando como base o “Modelo de Dados” da plataforma LA Referencia, validando as entidades, campos e relacionamentos aceitos pelo modelo.

3. RESULTADOS

Inicialmente, foram elencados os principais repositórios de dados que seriam fonte de informações para o BrCris. Diversos critérios foram ado-

tados para a seleção dos repositórios a serem utilizados, dentre eles: a consistência e atualização dos dados, o acesso aberto aos conjuntos de interesse, a amplitude dos repositórios e o reconhecimento dos dados pela comunidade científica brasileira. Como resultado, são agregados diversos repositórios, em diferentes formatos e com características distintas.

O principal repositório de dados para o BrCris são os currículos cadastrados na Plataforma Lattes do CNPq. Acessível em < <http://lattes.cnpq.br/> >, a Plataforma Lattes foi criada e é mantida pelo CNPq, contando atualmente com mais de 7 milhões de currículos cadastrados (em 15/04/2021), além de grupos de pesquisa e diretórios de instituições. Os currículos de interesse para o BrCris são aqueles dos pesquisadores, assistentes, técnicos que têm grau de mestre ou doutor, ou que possuem algum tipo de relação com a pesquisa científica (ou tecnológica), dado que possuem a publicação de artigo, compartilhamento de conjunto de dados, ou depósito de patente, ou têm alguma relação com a pós-graduação brasileira, seja como discente ou docente. São estimados, aproximadamente, 2,5 milhões de currículos de interesse.

De acordo com Lane (2010), em artigo publicado na revista Nature, a Plataforma Lattes é um poderoso exemplo de boas práticas para fornecimento de dados de alta qualidade. A autora relata também que órgãos federais, instituições e órgãos financiadores são incentivadores assíduos desta plataforma e que ela é uma das fontes de dados de pesquisadores mais confiáveis existente.

Além dos dados curriculares da Plataforma Lattes que subsidiaram informações para as entidades Person, OrgUnit, Publication, Patent, Event, Program, Course e Service, também são integrados dados dos seguintes repositórios:

- Oasisbr: mantido pelo IBICT, fornece dados confiáveis sobre publicações científicas em acesso aberto. Os dados foram mapeados para as entidades Publication, Service, Person.
- BDTD: a exemplo do Oasisbr, a BDTD também é mantida pelo IBICT. Fornece dados confiáveis sobre teses e dissertações brasileiras. Os dados foram mapeados para as entidades Publication, Course e Person.
- Plataforma Sucupira: concentra dados dos Programas de Pós-graduação do Brasil, fornecendo um conjunto de informações sobre os programas e cursos de pós-graduação. Todos os dados dos programas foram mapeados para as entidades Program, Course e OrgUnit;
- Instituições do INEP: assim como os programas de pós-graduação da Plataforma Sucupira, o INEP fornece uma base confiável, sobre as instituições de ensino do país em outros níveis de capacitação, sendo mapeados para a entidade OrgUnit.

- **Dados Abertos da CAPES:** fornece dados como publicações, orientações, entre outros que são mapeados para diversas entidades, como Person, OrgUnit, Publication, Program e Course.
- **Revistas Científicas:** o processamento de dados do conjunto de revistas científicas fornece informações diversas sobre elas, sendo mapeadas para a entidade Service. Exemplo de fonte de dados das revistas científicas são:
 - Diadorim
 - Latindex
 - DOAJ
 - UlrichsWeb

Como pode ser observado, as diversas fontes de dados mapeadas, se completam, possibilitando a criação de um conjunto padronizado e consistente, validado através de dados provenientes de diversas entidades brasileiras amplamente consolidadas e utilizadas. Ao se agregar todos os repositórios apresentados, é possível a adoção de técnicas que visam permitir a vinculação de conjuntos que inicialmente não era possíveis de se comunicarem, possibilitando dessa forma, a construção de um grande conjunto de dados, interligados, que facilitam a aplicação de consultas que inicialmente não seriam possíveis.

Diante do exposto, é possível verificar todo o conjunto de dados agregados no BrCris e também como estes foram selecionados e processados para uma integração, que engloba alguns elementos para possíveis desambiguações, utilizando para isso, identificadores implícitos ou gerados no processo de tratamento dos dados. Logo, com os dados agregados, diversas análises são viabilizadas, proporcionando um maior conjunto de análises bibliométricas.

Os resultados da execução do Projeto já incluem o desenvolvimento da arquitetura do BrCris, o mapeamento das fontes de dados a serem agregadas pelo Sistema, a implementação de provas de agregação dos recursos mapeados, a definição e realização de testes de serviços a serem disponibilizados. Entre as fontes agregadas, destacam-se em âmbito nacional, o Oasisbr, a BDTD, a Plataforma Lattes, a Plataforma Sucupira e o Portal de Dados Abertos da CAPES. Já entre as fontes internacionais, destaque é dado ao OpenAIRE Research Graph, DOIBoost, Portal Wikidata e ao DOAJ.

4. CONSIDERAÇÕES FINAIS

O BrCris se configura como um importante espaço de pesquisa e análise de dados. As informações agregadas e organizadas segundo um modelo

de dados semântico, permitem a geração de serviços para diversos atores, nos contextos de gestão e pesquisa acadêmica, assim como na área de informação para a inovação, que pretende ser o alvo da proposta apresentada. O BrCris é uma iniciativa que coleta e enriquece dados de repositórios e bases de dados abertas pela LA Referencia, utilizando protocolos OAI-PMH e múltiplos formatos de dados em XML e JSON. A próxima etapa do projeto é a aplicação dos sistemas de recomendações pelas métricas que podem ser explanadas em cada conjunto de dados.

CONFLITOS DE INTERESSE

Os autores declaram não haver conflito de interesses.

DECLARAÇÃO DE CONTRIBUIÇÃO

Administração do Projeto, Validação e Escrita: Washington Segundo.

Curadoria dos dados, Metodologia e Escrita: Thiago Dias, Tales Moreira.

Metodologia, Análise Formal e Investigação: Adilson Pinto, Vivian Silva, Josir Gomes, Luc Quoniam.

Análise Formal e Investigação: Lautaro Matas.

Curadoria dos dados: Ary Dias.

Administração do Projeto: Juliana Schneider

DECLARAÇÃO DE CONSENTIMENTO DE DADOS

O repositório de dados utilizados neste trabalho ainda não está disponível para compartilhamento. 

REFERÊNCIAS

- SIVERTSEN G. (2019). Developing Current Research Information Systems (CRIS) as Data Sources for Studies of Research. In: Glänzel W., Moed H. F., Schmoch U., Thelwall M. (eds) *Springer Handbook of Science and Technology Indicators*. Springer Handbooks. Springer, Cham. https://doi.org/10.1007/978-3-030-02511-3_25
- LANE, J. (2010). Let's make science metrics more scientific. *Nature*, 464(7288), 7288, p. 488-489. <https://doi.org/10.1038/464488a>

