

Repositórios de Dados: Objetivos, Funcionalidades e Alternativas

João Rocha da Silva*

joaorosilva@gmail.com

INESC TEC / Faculdade de Engenharia da Universidade do Porto
Porto, Portugal

RESUMO

Os repositórios de dados de investigação são cada vez mais uma peça essencial para o processo científico. Não só fomentam a reprodutibilidade das conclusões publicadas nos artigos científicos como assumem um papel crucial na atribuição de crédito aos criadores de dados, pois expõem ao público um trabalho de recolha, processamento e anotação que tanto tem de dispensioso como por vezes de invisível. A escolha de um software para suportar um repositório de dados deve ser guiada pelas necessidades das principais partes interessadas. Assim sendo, neste artigo discutem-se as principais funcionalidades desejáveis num repositório de dados, tanto do ponto de vista técnico (software e infraestrutura) como do ponto de vista político, nomeadamente no que diz respeito às garantias e compromissos a assumir pelas instituições que os alojam no sentido de permitir a sua certificação de acordo com a estratégia da European Open Science Cloud.

PALAVRAS-CHAVE

repositórios de dados, ciência aberta, e-Science, Estratégia para os Dados Abertos, Fundação para a Ciência e Tecnologia

1 INTRODUÇÃO

Os princípios Findable, Accessible, Interoperable, Reusable (FAIR) para a gestão de dados de investigação especificam que os dados devem ser *Findable* (encontráveis), *Accessible* (acessíveis), *Interoperable* (interoperáveis) e *Reusable* (reutilizáveis) [20]. Recentemente, um grupo de peritos da Comissão Europeia estudou as melhores práticas para tornar os dados de investigação europeus mais FAIR, tendo concluído que os componentes essenciais de um ecossistema FAIR são as políticas, os planos de gestão de dados (DMP), os identificadores persistentes, as normas, e finalmente os repositórios [6].

Sendo esta análise mais focada nos aspetos técnicos ligados às plataformas de software que podem sustentar um repositório, o aspeto político não será a tónica principal; ressalva-se contudo que as decisões políticas devem ser informadas pelas capacidades técnicas das plataformas disponíveis—é aliás nesse sentido que se este artigo se posiciona.

Falando de planos de gestão de dados, assiste-se a uma necessidade crescente de automação dos processos de gestão de dados, motivada pela grande quantidade e diversidade dos dados produzidos pelos investigadores. Os planos de gestão de dados consistem vulgarmente em documentos escritos e destinados a seres humanos, que são acrescentados às propostas de projetos de investigação em certos concursos. Mais recentemente, têm sido desenvolvidos esforços no sentido de os tornar acionáveis por máquinas, dando

origem aos machine-actionable Data Management Plan (maDMP). Estes apresentam-se como um documento vivo e interpretável por máquinas, reduzindo assim a necessidade de intervenção humana na execução e auditoria das práticas de gestão de dados a que os investigadores se comprometem aquando da criação do documento. Os repositórios devem portanto caminhar no sentido de suportar as operações necessárias à execução dos maDMP.

Os identificadores persistentes, ou Persistent Identifier (PID), são uma parte essencial dos metadados acrescentados aos conjuntos de dados, pois permitem a recuperação e identificação única de cada conjunto de dados. Qualquer plataforma de repositório deve então oferecer integrações com os fornecedores de PID. Há diversos sistemas de identificadores por onde escolher, e essa seleção deve também ser controlada por normas de certificação. A escolha da certificação prestigianete mas também acessível deve ser feita no sentido de aumentar o valor acrescentado para todas as necessidades e os recursos disponíveis das partes interessadas desse repositório. Devem portanto ser tidas em conta as raras mas preciosas histórias de sucesso na certificação aquando da seleção de uma plataforma de software para a montagem de um repositório.

Este artigo contém 5 secções para além desta introdução. A secção 2 apresenta uma análise de alto nível dos requisitos a satisfazer por um repositório de dados, tendo os princípios FAIR como linhas orientadoras. Tendo em conta a crescente importância da automação nos processos de gestão de dados como suporte aos Planos de Gestão de Dados accionáveis por máquinas, analisam-se também esses aspetos na secção 3. Na secção 4 é apresentada uma comparação entre soluções de software para repositórios. Esta foi produzida a partir de um estudo do grupo de trabalho RDA sobre plataformas de dados de investigação e uma comparação feita pela Dataverse, após consulta aos principais fornecedores de repositórios. Dada a importância da certificação no desenho de um *workflow* de gestão de dados, apresenta-se também na secção 5 uma discussão sobre a evolução de alguns dos esquemas de certificação de repositórios de dados, em linha com as mais recentes recomendações da Comissão Europeia. O artigo termina com a secção 6, onde se apresentam algumas conclusões e discussão sobre a análise levada a cabo, assim como algumas opiniões relativamente ao futuro dos repositórios de dados.

2 O PAPEL DO REPOSITÓRIO NA APLICAÇÃO DOS PRINCÍPIOS FAIR

Os princípios FAIR são linhas orientadoras para a melhoria dos processos de gestão de dados, abarcando portanto aspetos bastante mais genéricos do que a tecnologia que suporta um repositório de dados. Assim, e para ajudar a conduzir um artigo bastante focado

*Investigador no INESC TEC e Professor Auxiliar Convidado na Faculdade de Engenharia da Universidade do Porto

em tecnologia, utilizaremos os princípios FAIR como guiões orientadores para elencar algumas das funcionalidades relevantes de um repositório de dados.

O primeiro dos princípios FAIR diz que os dados devem ser *Findable*, ou seja, devem ser fáceis de descobrir, tanto por humanos como por máquinas. Para tal, devem incluir (F1) um identificador global único e persistente, (F2) possuir metadados ricos, que por sua vez (F3) devem incluir claramente o identificador dos recursos que descrevem. O quarto e último aspeto, (F4), especifica que os tanto dados como metadados devem estar registados em recursos pesquisáveis.

Para facilitar a satisfação destes requisitos, uma plataforma de repositório deverá oferecer integração programática com um fornecedor de identificadores persistentes como por exemplo os Digital Object Identifier (DOI) ou handle¹. A vantagem para os gestores de repositório é a possibilidade de atribuir identificadores automaticamente aos conjunto de dados disponibilizados. Um identificador persistente funciona, em palavras simples, como um atalho para um programa num computador. Não é uma cópia do conjunto de dados, mas sim apenas um apontador para o local onde esse conjunto de dados está publicado (tipicamente, um repositório). Sendo assim, deve poder ser *desreferenciado*, ou seja, deve ser possível obter o conjunto de dados a partir do seu identificador. Para um ser humano, tal operação pode consistir apenas num simples clique numa ligação apresentada no seu navegador Web; para uma máquina, contudo, pode ser um pedido de rede à plataforma que gere os identificadores. Em ambos os casos, a entidade solicita o recurso por detrás daquele identificador; no caso do humano, a resposta será em HyperText Markup Language (HTML), código que o seu navegador pode interpretar de forma a construir uma página web legível para humanos. Para uma máquina, o formato devolvido poderá ser, por exemplo eXtensible Markup Language (XML) ou Resource Description Framework (RDF) (formatos de representação de informação que as máquinas podem interpretar).

No caso do repositório de dados do INESC TEC², a Fundação para a Ciência e Tecnologia (FCT) negociou um contrato para a aquisição de pacotes de identificadores em colaboração com a DataCite³, que permitiu ao repositório emitir DOI para os conjuntos de dados depositados.

Para satisfazer os requisitos F2 e F3 é necessário que os registos de metadados associados a cada conjunto de dados obedeçam a um esquema normalizado, como é o caso do DataCite Schema⁴. Uma grande vantagem do uso de um esquema de metadados em conjugação com a atribuição de um DOI é que cada identificador ficará associado a uma ficha de metadados no momento da cunhagem do DOI. Tipicamente, as plataformas de emissão de identificadores não permitem sequer a sua cunhagem sem o preenchimento dessa ficha de metadados. Esta ficha ficará então guardada na plataforma emissora do DOI, e mesmo que o conjunto de dados deixe de estar disponível no repositório, a ficha continuará disponível para consulta na plataforma fornecedora do identificador (satisfazendo o requisito A2).

O segundo princípio FAIR impõe que os dados sejam *Accessible*, pois tanto dados como metadados devem ser acessíveis através de identificadores, através de protocolos de comunicações normalizados (A1). Este ponto requer que o protocolo seja aberto, livre e universalmente implementável (A1.1), e que inclua procedimentos de autenticação e autorização, quando necessário (A1.2). O segundo requisito para assegurar a acessibilidade (A2) é que os metadados devem permanecer acessíveis, mesmo quando os dados deixam de o estar.

Para satisfazer o requisito A1, um repositório deverá suportar o processo de desreferenciação de identificadores persistentes, quer por máquinas quer por seres humanos. Quando se fala de protocolos de comunicações normalizados, fala-se quase sempre do protocolo HyperText Transfer Protocol (HTTP). Este é talvez o protocolo mais usado na web para transferência de informação, e satisfaz o requisito A1. Ele inclui dois mecanismos relevantes para a satisfação do requisito A1.1, chamados *Content Negotiation*⁵, *Authorization*⁶ e *Access Authentication*⁷.

De uma forma muito breve, este mecanismo permite a um cliente (pode ser um navegador web ou um programa de computador) solicitar ao servidor (máquina que disponibiliza o registo de metadados de um determinado conjunto de dados) que envie a informação determinado formato, ao escrever esse pedido no cabeçalho do pedido HTTP (que pode ser visto como a secção “Destinatário” de um envelope). Ao especificar que pretende HTML, um navegador irá obter o código necessário para apresentar os dados sobre o conjunto de dados a um ser humano; se uma máquina solicitar XML ou RDF, irá obter um documento que é muito mais difícil de ler por parte de um ser humano, mas que uma máquina interpretará corretamente. A beleza deste sistema de negociação é que permite obter essas diferentes representações a partir de um mesmo identificador, mudando apenas o cabeçalho do pedido inicial.

Apesar da autenticação básica oferecida pelo protocolo HTTP ser suficiente para satisfazer requisito A1.1, existem protocolos alternativos que podem ser adoptados caso hajam outros requisitos relevantes para a instituição que implementa um repositório de dados. Listam-se alguns protocolos abertos que podem ser relevantes neste cenário:

- **Autenticação federada:** O protocolo Shibboleth⁸ é bastante usado nas instituições de ensino superior portuguesas para controlar o acesso a múltiplos recursos através do portal de autenticação federada oferecido pela Fundação para a Computação Científica Nacional (FCCN). Um repositório institucional deverá registar-se junto do Identity Provider relevante para permitir aos utilizadores autenticar-se com as mesmas credenciais que usam para aceder aos restantes recursos da sua instituição.
- **Autenticação delegada via outros providers:** O protocolo OAuth 2.0⁹ permite aos utilizadores autenticar-se com as mesmas credenciais que utilizam em outros serviços como

¹<http://www.handle.net/>

²Ligação: <https://rdm.inesctec.pt>

³Ligação: <https://datacite.org/>

⁴Ligação: <https://schema.datacite.org/meta/kernel-4.3/>

⁵Ligação: <https://www.w3.org/Protocols/rfc2616/rfc2616-sec12.html>

⁶Ligação: <https://www.w3.org/Protocols/HTTP/1.0/draft-ietf-http-spec.html#Authorization>

⁷Ligação: <https://www.w3.org/Protocols/HTTP/1.0/draft-ietf-http-spec.html#BasicAA>

⁸Ligação: <https://www.shibboleth.net/index/>

⁹Ligação: <https://oauth.net/2/>

o Open Researcher and Contributor ID (ORCID)¹⁰, que pode ser útil para assegurar a autenticidade de auto-depósitos ou modificações feitas aos metadados pelos próprios criadores dos mesmos. Desta forma, podemos garantir que se uma determinada operação é realizada por alguém, se esse alguém se tiver autenticado com sucesso junto do ORCID, por exemplo. Este mecanismo foi implementado na plataforma de gestão de dados Dendro, para agilizar o registo de novos utilizadores e a sua autenticação. Desta forma, o ORCID de cada utilizador ficará associado ao seu perfil, e indiretamente às operações efetuadas por esse utilizador no repositório.

O terceiro princípio FAIR diz que a gestão de dados deve ser *Interoperable*. Isto deve-se ao facto dos dados precisarem frequentemente de ser integrados com outros dados, e como tal têm que ser facilmente integráveis em fluxos de trabalho de processamento. Só assim podem ser facilmente armazenados e analisados de forma interoperável. Para assegurar a interoperabilidade, tanto dados como metadados devem (I1) usar uma linguagem formal, acessível e largamente aplicada para representação de conhecimento, (I2) devem usar vocabulários que, por sua vez, seguem os princípios FAIR, e (I3) devem incluir referências qualificadas para outros dados e metadados.

Talvez o princípio mais difícil de assegurar, a interoperabilidade implica a representação da informação constante no repositório em formatos que permitam a interpretação por sistemas externos. Isto quer dizer que os dados devem ser representados não só em formatos amigáveis como largamente suportados pelas bibliotecas de manipulação de dados mais usadas e disponíveis em código aberto. Para os dados, por exemplo, devem ser adotados os formatos livres de dependências e largamente suportados, como por exemplo XML, Comma-separated Values (CSV) ou Tab-separated values (TSV), em detrimento de formatos binários ou proprietários. Para os metadados, devem ser representados com recurso a standards bem definidos na comunidade, sejam eles formalizados como esquemas XML, ou ontologias no caso do repositório disponibilizar os registos de metadados como Linked Open Data (LOD).

Alguns exemplos destes esquemas incluem o amplamente utilizado Dublin Core ou o schema.org, um esquema de metadados para a web que reuniu a colaboração das principais tecnológicas (Google, Microsoft, Yahoo e Yandex)¹¹, que possam tanto ser formalizados como esquema XML como sob a forma de ontologias. Sejam qual forem os vocabulários seleccionados, eles próprios devem seguir os princípios FAIR, o que significa que devem estar livremente acessíveis e devidamente documentados.

Por último, a interoperabilidade também diz respeito à inclusão de referências para recursos relacionados. Dessa maneira, deverão ser incluídas nos registos de metadados de cada conjunto de dados referências para materiais relacionados. Exemplos de materiais relacionados incluem por exemplo um documento de dissertação ou tese, artigos resultantes da produção dos conjuntos em questão. Os esquemas ou ontologias utilizadas para a descrição devem assim incluir descritores que permitam essas referências, como é o

exemplo do descritor *references* do Dublin Core¹² ou o descritor *citation* do schema.org¹³.

O quarto e último dos princípios defende que os dados devem ser *Reusable*, pois o objetivo final dos princípios FAIR é fomentar a sua reutilização. Para serem Reutilizáveis, tanto dados como metadados devem (R1) ser descritos com uma grande variedade de atributos corretos e relevantes. Este princípio subdivide-se em três: (R1.1) devem ser publicados com uma licença de utilização clara e acessível, (R1.2) devem estar associados a informação de proveniência detalhada, sendo que esses dados e metadados devem satisfazer as normas relevantes de cada domínio (R1.3).

As licenças a escolher (recomendação R1.1) ficam ao critério da instituição de investigação, da entidade financiadora, e de quem cria os dados. Contudo, para a disponibilização de dados de investigação em regime de *Open Data* é desejável a adoção de licenças que permitam a reutilização dos dados e metadados com o mínimo de restrições.

Para oferecer a terceiros a liberdade de reutilizar os dados em trabalhos derivados e ao mesmo tempo assegurar a atribuição de crédito aos autores, as licenças Creative Commons (nomeadamente a CC BY-4.0)¹⁴ são geralmente adequadas para a disponibilização de conjuntos de dados produzidos em projetos financiados por dinheiros públicos. Grande parte dos investigadores podem, no entanto, sentir dificuldades no momento da escolha da licença mais adequada para os conjuntos de dados que desejam publicar. É nesta altura que uma plataforma de repositório pode complementar o trabalho dos curadores, ao fornecer uma lista das licenças mais comuns incluindo descrições claras e sucintas de cada licença. Para aqueles utilizadores que pretenderem uma análise mais pormenorizada devem também oferecer uma ligação direta para a licença completa em linguagem jurídico-legal.

A recomendação R1.2 refere a necessidade de manter um historial que permita estabelecer a proveniência de um conjunto de dados. Enquanto que é relativamente simples manter um historial de modificações dentro de uma plataforma de repositório, o desafio é maior quando se considera que um conjunto de dados pode ser derivado de outro, tendo este último o seu próprio historial de modificações registado em outra plataforma.

Existe aqui uma clara vantagem em usar LOD como formato de exposição de metadados (incluindo a proveniência) neste cenário. Para tal, é necessário que os repositórios de dados ofereçam um histórico de modificações em LOD que obedeça a ontologias como a PROV-O[10] (mais normativa pois é uma recomendação da World Wide Web Consortium (W3C)) ou a Provenance And Versioning (PAV) [3] (uma alternativa mais simples, também apelidada de *lightweight ontology*). Desta forma, quando uma terceira máquina des-referenciar o conjunto de dados derivado será também interrogado o repositório onde se encontra o conjunto de dados original, e assim sucessivamente, se ele próprio também for um conjunto de dados derivado. Torna-se assim transparente e automática a recuperação de todos os dados relativos ao historial de modificações e versões de um conjunto de dados, sem necessidade de indexar todos os recursos e sem que haja lugar à intervenção humana.

¹⁰Ligação: <https://members.orcid.org/api/oauth2>

¹¹Ligação: <https://schema.org/>

¹²Ligação: <https://www.dublincore.org/specifications/dublin-core/dcmi-terms/terms/references/>

¹³<https://schema.org/citation>

¹⁴Ligação: <https://creativecommons.org/licenses/by/4.0/>

É bastante claro que domínios de investigação distintos poderão gerar conjuntos de dados muito diferentes também. Como tal, o processo de descrição dos mesmos deve também ser adaptado. Satisfazer a recomendação R1.3, que preconiza que os dados e metadados devem obedecer às normas vigentes em cada domínio requer suporte por parte do software de repositório, que deve ser flexível ao ponto de permitir a parametrização de múltiplos esquemas de metadados para distintos domínios. Deverá também permitir a escolha e combinação de múltiplos descritores (genéricos e específicos do domínio do conjunto de dados em depósito) aquando do preenchimento do registo de metadados, sejam esses descritores genéricos ou específicos de domínio. É preciso também não esquecer o trabalho fundamental dos curadores no apoio institucional aos seus investigadores, pois a escolha dos descritores mais adequados a cada conjunto de dados.

3 A AUTOMAÇÃO DOS PLANOS DE GESTÃO DE DADOS

Qualquer processo de gestão de dados de investigação deve começar com o desenho de um Data Management Plan (DMP). Um DMP é um documento tipicamente escrito e aprovado antes do início de um projeto de investigação. Em alguns casos, é já requisito obrigatório em chamadas de propostas de projetos de investigação suportados por fundos públicos [16]) e que especifica todos os aspetos relacionados com a gestão dos dados produzidos ao longo do projeto [14].

Entre outros aspetos, o DMP especifica quais os dados a produzir no contexto de um projeto, que metadados serão anexados aos dados publicados, onde eles irão ficar disponíveis após o fim do projeto, entre outros pontos. No Reino Unido, o Digital Curation Centre (DCC) disponibiliza tanto um guião [8] para a escrita deste tipo de documentos como uma ferramenta online, o DMPOnline¹⁵, que assiste os investigadores e curadores na elaboração dos seus DMP, assistidos por uma base de dados de templates de DMP aceites por diferentes entidades financiadoras.

Com a grande quantidade de dados de investigação gerada diariamente surge a necessidade de acompanhar e monitorizar a implementação de DMP, mas de forma automática. Três exemplos simples de operações frequentes, repetitivas e facilmente automatizáveis são:

- Simples verificação da disponibilidade dos dados depositados
- Controlo de qualidade dos metadados associados aos conjuntos de dados acabados
- Verificação das licenças associadas aos conjuntos de dados

Para lidar com estas operações e outras de forma mais automática surgiram os chamados DMP accionáveis por máquinas (maDMP). O objetivo destes documentos é libertar os curadores de dados para outras operações menos automatizáveis, como o suporte direto aos investigadores ou atividades de formação para ajudar à curadoria de dados no dia-a-dia.

Como preconizado pelos autores das 10 regras para maDMP [15], deve ser possível aos sistemas informáticos levar a cabo ações em nome das partes interessadas no processo de gestão de dados de

investigação. Operações dadas como exemplo incluem a recolha automática de informação administrativa relevante para anotação dos recursos produzidos. Essa informação inclui a recuperação de referências aos financiadores dos projetos, os currículos dos autores dos conjuntos de dados ou as informações corretas das instituições de acolhimento dos investigadores.

Esta necessidade de suporte à automação e integração tem implicações vastas do ponto de vista do desenho e funcionalidades de um repositório de dados, na medida em que a solução deve disponibilizar uma Application Programming Interface (API) completa e bem documentada, de forma a permitir a sistemas externos executar estas operações sem a necessidade de intervenção humana.

3.1 O papel da interoperabilidade na descoberta de conjuntos de dados

Os agregadores de repositórios são portais que indexam o conteúdo dos repositórios para facilitar a pesquisa e descoberta de recursos neles contidos. Exemplos destes agregadores são, por exemplo, o portal re3data da DataCite¹⁶ ou o Dataset Search da Google¹⁷. Para facilitar a integração com diretórios de dados, o software de repositórios devem ser assentes em protocolos de interoperabilidade standard.

Para facilitar a indexação dos conteúdos dos repositórios, estes implementam suporte ao protocolo OAI-PMH (Open Access Initiative Protocol for Metadata Harvesting). Este protocolo expõe os metadados de todos os registos no repositório de forma paginada, para que seja lida sequencialmente e indexada; contudo, não possibilita a pesquisa e recuperação de registos específicos por termos contidos nos registos de metadados, por exemplo. Para conseguir tal funcionalidade é necessário proceder à inserção de todos os documentos num índice de pesquisa como o Lucene ou Solr. Uma consequência óbvia é a necessidade de manutenção de cópias dos registos e atualizações periódicas de todo o índice—uma operação lenta e dispendiosa em termos de recursos informáticos, que coloca todo o custo de manutenção do serviço do lado dos clientes.

Uma situação oposta ocorre com a adoção de LOD na representação de metadados para conjuntos de dados de investigação. Esta representação de dados vem eliminar a necessidade de indexação periódica de conteúdos e melhorar muito a precisão das pesquisas, caso os dados sejam estruturados de acordo com uma ou mais ontologias. Uma ontologia é definida, na Ciência da Informação, como uma *especificação de uma conceptualização*[4], que permite às máquinas interpretar o significado das entidades representadas num sistema e das relações existentes entre elas.

Uma representação via LOD permite uma maior descentralização da carga de pesquisas sobre os servidores onde os repositórios estão alojados. Reduz também a necessidade de indexação porque é possível executar pesquisas nos próprios servidores sem ter que primeiro indexar todo o seu conteúdo e permite aos clientes do serviço recuperar conjuntos de dados com muito mais precisão do que nos casos em que se usa um índice de pesquisa simples baseado em palavras-chave, pois é possível especificar critérios muito detalhados sob a forma interrogações SPARQL Protocol and RDF Query Language (SPARQL).

¹⁵Ligação: <https://dmponline.dcc.ac.uk>

¹⁶Ligação: <https://www.re3data.org/>

¹⁷Ligação: <https://datasetsearch.research.google.com/>

A disponibilização de LOD por parte dos repositórios tem, por contraponto aos benefícios de interoperabilidade e abertura, um custo para quem os hospeda. As consultas SPARQL são exigentes do ponto de vista computacional e, se um número elevado de utilizadores (máquinas ou humanos) tentarem ao website ao mesmo tempo, rapidamente aumentarão os tempos de resposta da infraestrutura. Isto é especialmente relevante quando se considera a funcionalidade de pesquisa federada, onde múltiplos servidores podem ser contactados para conseguir obter os resultados pretendidos por uma única pesquisa iniciada em qualquer dos nós da federação.

O futuro dos repositórios de dados pode passar então por uma solução híbrida, onde servidores e clientes partilham informação para possibilitar também uma partilha do esforço de interrogação e descoberta de dados. Soluções como a Linked Data Fragments[18] propõem a combinação de interrogações sobre SPARQL *endpoints* tradicionais (no servidor do repositório que expõe LOD) com a interrogação a dados locais a cada cliente, reduzindo assim a carga sobre os servidores.

3.2 Repositórios como plataformas de computação

Cada vez mais, o valor de um repositório de dados reside na facilidade de reutilização dos dados nele contidos. Tecnologias recentes como os Research Notebooks e a Containerização permitem tornar os algoritmos mais portáteis, de forma a poderem ser executados remotamente sobre os dados depositados num repositório. O papel do repositório torna-se assim muito mais do que um armazém combinado com um motor de busca sobre conjuntos de dados.

Os Research Notebooks são cadernos de laboratório electrónicos que combina num único pacote os dados-base e o código de processamento desses dados. Esse código pode servir, por exemplo, para transformar os dados base em dados processados, e até mesmo para produzir visualizações e gráficos em tempo real [13].

Os Notebooks vêm também possibilitar suportar um novo tipo de publicação científica, fortemente suportada por dados e pela Internet. A combinação elegante de texto, trechos de código e visualizações apelativas permitiu a criação dos chamados “Web Journals” como por exemplo o *distill.pub*¹⁸. Estes distinguem-se das publicações convencionais ao embutir visualizações interativas no texto dos artigos científicos, fomentando a experimentação através da manipulação de parâmetros de entrada dos algoritmos que geram essas visualizações, para produzir um feedback visual imediato aos leitores. Toda a computação tem que ser executada em tempo real, o que implica a montagem de uma réplica do ambiente de processamento por detrás do portal web que suporta o journal.

Com o empacotamento do processamento e dos dados a ser abordado pelos Research Notebooks, surge a necessidade de otimizar a localização da computação e dos dados, para aproximar um do outro. Nem todos os conjuntos de dados se prestam a ser transportados pela rede em tempo útil para a execução de um determinado algoritmo, pois em certas disciplinas o seu tamanho pode ascender a centenas de giga ou terabytes.

Torna-se necessário levar a computação até aos dados em vez de ter que transmitir os dados até ao local onde se realiza a computação. Este tipo de computação que ocorre junto dos dados é

um paradigma é implementado em portais de computação como o D4Science¹⁹ ou o EUDAT B2Stage²⁰. Com a esperada descentralização e interoperabilidade que se espera dos repositórios de dados defendida pelo Guia Estratégico de Implementação (GEI) da European Open Science Cloud (EOSC), espera-se que cada repositório seja capaz de oferecer capacidade de computação local a pequena escala e junto das fontes de dados. Este cenário é mais próximo do *edge computing*, um termo mais vulgarmente associado às aplicações Internet of Things (IoT). Por contraponto ao *cloud computing*, onde todos os recursos se encontram na nuvem e os dados têm que ser cuidadosamente aproximados dos nós de computação (mesmo em termos geográficos) para reduzir tráfego na rede, o *edge computing* propõe que cada nó da rede deve ter uma pequena mas importante capacidade de computação para executar operações à medida que os dados são produzidos ou atualizados.

Quando se fala de distribuição de computação e reprodutibilidade, surge a necessidade de replicar também todo o ambiente de execução no qual essa computação é executada. Sem uma solução de virtualização, isso implica a instalação do sistema operativo e de todas as dependências que foram originalmente usadas na execução do código. Este é um processo manual e demasiado dispendioso, criando assim a necessidade de assegurar a portabilidade do ambiente de execução para o tornar automaticamente instalável em qualquer nó.

A containerização é uma tecnologia de virtualização que, quando comparada com as máquinas virtuais convencionais, elimina a necessidade de virtualizar um sistema operativo para executar um determinado programa, substituindo essas dependências por um ambiente de execução (ou *runtime*) próprio. O exemplo mais comum de ambiente de virtualização é o Docker²¹. Esta eficiência torna possível a um repositório de dados executar algoritmos directamente sobre os dados nele depositados e retornar apenas os resultados aos clientes que pedem esses processamentos, em vez de exigir a descarga dos dados para o cliente e colocar sobre este a responsabilidade de montar todo o ambiente de execução.

Como parte da especificação do ambiente de containers, é necessário de especificar todos os passos de instalação de dependências necessárias à execução dos algoritmos de processamento de dados. Este código—que acaba por ser uma espécie de documentação funcional—reduz em larga medida os problemas de replicação do ambiente de execução de determinado algoritmo.

3.3 A blockchain e as suas aplicações aos Dados Abertos

Num futuro de repositórios de dados interligados num ambiente descentralizado será necessário conceber soluções para assegurar a autenticidade da informação colocada em cada nó dessa rede. Em particular, a criação de um *ledger* de transações distribuído—como uma *blockchain* envolvendo os repositórios na rede—pode reduzir a dependência de entidades terceiras—como os emissores de certificados—para assegurar a autenticidade dos dados. Este tema

¹⁸Ligação:

¹⁹Ligação: <https://www.d4science.org>

²⁰Ligação: <https://www.eudat.eu/b2stage>

²¹Ligação: <https://www.docker.com/>

já está a ser alvo de investigação, com o objetivo de descobrir formas de utilização da blockchain para assegurar a validade dos resultados de operações sobre dados abertos [17]. É portanto claro o potencial de aplicação deste tipo de tecnologias aos diversos aspetos da gestão de dados abertos, destacando-se dois pontos:

- **Proveniência:** Através de um registo distribuído de modificações feitas a dados e metadados será possível assegurar a sua autenticidade e não-repúdio relativamente à autoria de quaisquer operações realizadas sobre os dados. Esta prática deve estar presente antes, durante e após a publicação dos conjuntos de dados; quando aplicada antes da publicação, permite rastrear todas as modificações feitas ao longo do tempo, para reduzir situações de *p-hacking*[5], por exemplo. Após a publicação, um registo distribuído de alterações permite aos autores continuar a modificar um conjunto de dados após a publicação do conjunto de dados considerado “final”, de forma completamente transparente. Uma outra vantagem é a resiliência do sistema, pois mesmo que o repositório onde um conjunto de dados está publicado seja desligado, esse historial de modificações continuará registado nos restantes nós da *blockchain*.
- **Fiabilidade de resultados:** Regra geral, um valor só é registado na *blockchain* se existir um consenso na rede de computadores que colaboram nessa *blockchain*. Este mecanismo pode ser usado para atestar a reprodutibilidade das operações de transformação realizadas sobre um conjunto de dados. Neste caso em particular, só após um determinado número de nós (dependendo da política em vigor na *blockchain*) levarem a cabo as mesmas operações e confirmarem os resultados é que haverá o *consenso* necessário à escrita para a *blockchain*, passando aquela operação a fazer parte do historial de operações sobre os dados envolvidos.

4 COMPARAÇÃO DE FUNCIONALIDADES

Existem diferentes alternativas para a montagem de um repositório de dados. Uma análise comparativa produzida pelo projeto DataVerse²² apresenta uma comparação das principais plataformas de repositórios de dados, tendo em conta três grupos de facetas: Funcionalidades de Software, Controlo/Organização e Conteúdo. Uma outra comparação [1] discute em mais profundidade algumas características técnicas de diversas plataformas. Esta última foca-se apenas nas alternativas disponíveis em regime de código aberto e instaláveis localmente, de forma a poder retirar conclusões relativamente ao modelo de dados por detrás de cada solução de software.

Tão importante como comparar funcionalidades oferecidas pelas plataformas é estudar quais funcionalidades são mais valorizadas pelos seus utilizadores. Para responder a esta questão, o Repository Platforms for Research Data Interest Group (RPRD IG) da Research Data Alliance (RDA) compilou uma Matriz de Casos de Uso e Requisitos Funcionais para as Plataformas de Gestão de Dados de Investigação. Esta matriz resultou de um inquérito a 11 grupos que apresentaram os seus casos de uso para repositórios de dados de investigação, tendo cada um classificado a importância de 26 requisitos vulgarmente associados a um repositório de dados. No final, foi atribuída uma pontuação a cada um dos requisitos funcionais,

com base na importância que a totalidade dos grupos lhe atribuiu no contexto do seu caso de uso.

De forma a comparar de forma sucinta algumas das plataformas mais utilizadas, é necessário considerar não só a quantidade e variedade de funcionalidades como também a sua importância para a maioria dos utilizadores. Assim, na Figura ?? apresenta-se não só uma pontuação de funcionalidades, como também algumas das principais características dessas plataformas. Nesta figura apresenta-se primeiramente um gráfico de barras que ordena as plataformas de acordo com uma “Pontuação de funcionalidades”. Esta foi determinada associando manualmente os requisitos funcionais apresentados pelo RPRD IG às características listadas na comparação de repositórios apresentada pela DataVerse, e calculando o produto de cada relação estabelecida pela importância dada ao requisito pela RDA. Desta forma, foi possível atribuir diferentes importâncias ou pesos às características desta última comparação²³.

A figura apresenta também uma tabela com algumas características mais relevantes de cada plataforma:

• Código Aberto

Na primeira linha distinguem-se as soluções que estão disponíveis em regime de código aberto daquelas em que não é o caso. O facto do código de uma solução estar disponível em acesso aberto tem inúmeras vantagens, desde que haja uma comunidade disposta a manter esse software. Exemplos são uma maior segurança, pois é possível a qualquer utilizador escrever e publicar correções para falhas de segurança, ao invés de ter que esperar pela correção por parte da equipa de desenvolvimento de uma empresa. Qualquer pessoa tem também a liberdade de modificar e melhorar o software para satisfazer os seus requisitos específicos. Esses melhoramentos podem até ser incluídos no tronco comum da solução, caso se verifique que o requisito está presente em outras comunidades e os desenvolvedores assim o entenderem. Esta flexibilidade pode proporcionar um ritmo e velocidade de desenvolvimento mais elevado do que as soluções proprietárias, que têm equipas de desenvolvimento limitadas, mas a comunidade torna-se responsável pela manutenção da qualidade do código. O modelo de código aberto não invalida a hipótese das entidades que suportam um repositório realizarem alterações pagas, num cenário semelhante à consultoria de implementação do repositório. Nestes casos, o “produto” da empresa em si não é o software, mas sim o serviço de instalação e personalização da solução de código aberto numa determinada instituição. O código desenvolvido no contexto destas implementações pagas pode até ser incluído no tronco comum publicamente disponível, beneficiando assim toda a comunidade que utiliza esse software. É também normal encontrar nos projetos de software de código aberto um conjunto de fóruns públicos onde os utilizadores podem propor melhorias ao software e reportar *bugs*.

- **Grátis** Quase todas as soluções em análise são completamente grátis ou oferecem um subconjunto de funcionalidades grátis. Os sistemas que funcionam com base na *nuvem*

²²Ligação: <https://dataverse.org/blog/comparative-review-various-data-repositories>

²³Os dados-base deste estudo estão disponíveis no repositório GitHub deste relatório, como materiais auxiliares. Ligação: <https://github.com/silvae86/repositories-paper-eda>

Figura 1: Análise comparativa de diferentes plataformas de *software* para repositórios de dados

incluem normalmente uma capacidade de armazenamento de dados limitada, sendo que os utilizadores podem adquirir um pacote de armazenamento maior se necessitarem. As soluções que podem ser instaladas localmente apenas dependem da capacidade de armazenamento do servidor no qual são instaladas.

- **Número de conjuntos de dados, ficheiros e utilizadores** A quantidade de recursos geridos pelas plataformas e o número de utilizadores que confiam nelas para a sua gestão de dados é uma importante métrica da sua maturidade.
- **Certificação** A certificação de repositórios continuará a assumir relevância aquando da adoção de uma solução de repositório. No caso das soluções Software as a Service (SaaS), a certificação do repositório pode beneficiar todas as instituições que contratem os serviços dessas plataformas, pois todos os seus investigadores poderão mencionar esse aspeto na escrita dos DMP que anexam às suas propostas de projetos. As soluções SaaS têm assim interesse em procurar a certificação para oferecer mais valor aos seus clientes. No caso das alternativas instaladas e mantidas localmente, as funcionalidades da solução de software em particular não são condição suficiente para a certificação, mas as funcionalidades

da plataforma a seleccionar (ou falta delas) não devem ser um obstáculo a futura certificação do repositório.

- **Infraestrutura** Esta categoria distingue as plataformas que só estão disponíveis como SaaS daquelas que podem também ser instaladas e geridas localmente por uma equipa interna à organização e suportada em recursos computacionais locais.

Por fim, é necessário considerar o modelo de hospedagem ou infraestrutura do repositório, que está geralmente ligado ao modelo de negócio das entidades que o suportam. Neste domínio, existem as soluções SaaS, que podem ser vistos como “chave-na-mão”. Nestes modelos, toda a infraestrutura e serviços por detrás da página web do repositório (alojamento, manutenção, assistência técnica, curadoria ou atribuição de identificadores) faz parte de um pacote oferecido às instituições interessadas nesse repositório. Em alternativa, apresentam-se as soluções alojadas localmente, que diferem das soluções SaaS na medida em que exigem instalação e manutenção por parte de uma equipa de TI interna à organização. Da mesma forma, o trabalho de curadoria na plataforma também tem que ser levado a cabo inteiramente pelos elementos da organização e geralmente não existe nenhuma obrigação de assistência por parte dos desenvolvedores do software.

Tabela 1: Recursos dispendidos na certificação do TIB, para o DSA V.2 e nestor Seal, números [12]

Métrica	DSA	nestor Seal
Fee for process	0 €	500€
Project Duration (meses)	9	12
Pessoas-Mês ²⁴	3.7	11
Pessoas envolvidas	7	16
Unidades organizacionais envolvidas	5	8

5 CERTIFICAÇÃO

Para além dos princípios FAIR, que se apresentam como linhas gerais para aquilo que um processo de gestão de dados deve proporcionar, surgiram recentemente os princípios TRUST [11]. Estes focam-se no problema da confiança nas infraestruturas, organizações e qualidade dos dados de investigação. Estes princípios oferecem, de acordo com os autores, uma *framework* de discussão para ajudar a melhorar a confiança de todas as partes interessadas no processo de gestão de dados. Estas partes interessadas vão desde os próprios investigadores, passando pelas instituições académicas ou de ciência que os acolhem, as entidades financiadoras, entre outros [19].

De acordo com o Guião Estratégico de Implementação (GEI) da EOSC (European Open Science Cloud), a certificação de repositórios de dados é uma parte essencial na visão a longo prazo da Comissão Europeia. Os processos de certificação são exaustivos e requerem não só excelência técnica como também um compromisso a longo prazo das instituições que os suportam.

De um ponto de vista mais económico e de soberania dos dados, a certificação de repositórios assumirá um papel cada vez mais importante no funcionamento de uma instituição académica ou de investigação. Existem diversas razões, de entre as quais se destaca o crescente número de entidades financiadoras e publicações de renome que exigem a disponibilização dos dados-base e do código-fonte que sustenta as publicações. Neste cenário, a qualidade dos repositórios onde tais recursos são disponibilizados terá necessariamente que ser elevada. Desta forma, a certificação tornar-se-á assim uma pedra fundamental no suporte aos projetos de investigação, pois as instituições terão que assegurar aos seus investigadores a satisfação dos presentes e futuros requisitos de financiadores aquando da escrita dos planos de gestão de dados a anexar às propostas de financiamento. Essas instituições terão então duas opções para garantir os requisitos: ou embarcam num processo de certificação dos seus próprios repositórios institucionais, ou terão que subcontratar o processo de depósito e curadoria de dados a um repositório certificado. Levada a cabo a certificação, contudo, o repositório poderá tornar-se uma fonte de receita para as instituições que os suportam, pois oferece um valor acrescentado caso essas instituições queiram oferecer um serviço de preservação a longo prazo como um serviço [12].

A certificação é um processo custoso e demorado que constitui um projeto por si só. A título de exemplo, o TIB Leibniz Information Centre for Science and Technology, que reportou os resultados das suas certificações DSA e nestor em 2018; os números são apresentados na Tabela 1.

5.1 Data Seal of Approval

O Data Seal of Approval (DSA)²⁵ foi uma iniciativa de certificação de repositórios desenvolvida em 2008 pelo holandês Data Archiving and Networked Services (DANS). Consiste num selo que garante que dados arquivados num repositório DSA podem ser encontrados, interpretados e usados no futuro. De 2008 a 2018, os requisitos para obtenção do selo passaram por 3 revisões. O DSA foi fundido em 2018 com o CoreTrustSeal, e os repositórios que receberam o DSA para o período 2014-2017 têm que passar novamente pelo processo de certificação para o CoreTrustSeal. Uma lista dos repositórios que obtiveram o DSA durante a sua existência está também disponível²⁶.

5.2 nestorSEAL

A nestor (network of expertise in long-term storage of digital resources in Germany) envolve 23 grandes instituições alemãs de investigação e ciência e deu origem a 12 grupos de trabalho. Um grupo de trabalho especializado em certificação mantém o nestor-SEAL, uma *framework* de certificação que garante não só os requisitos do DSA mas também uma certificação mais extensa²⁷.

5.3 ISO 16363:2012

A certificação ISO 16363:2012 assume-se como a mais exigente e valiosa, mas também a dispendiosa e demorada de obter. Nem todos os repositórios devem procurar esta certificação, pois como diz o Implementation Guide da EOSC, “um repositório deve procurar o nível de certificação apropriado e alcançável”[6]. A certificação ISO é um processo considerado “heavyweight” para a maioria dos repositórios, e até 2018 apenas o National Cultural Audiovisual Archives, na Índia, havia atingido esse objetivo [7, 12].

Uma auditoria interna nacional realizada em 2015 analisou 24 repositórios alojados pelo SARI²⁸ do RCAAP²⁹, de acordo com as três dimensões da norma: infraestrutura organizacional, gestão de objetos e infraestrutura de gestão e segurança.

Durante esta auditoria, e dado que a ISO 16236:2012 não especifica uma escala de maturidade para o cumprimento dos requisitos nela especificados, foi adotada uma escala baseada no modelo ECM3 [9]. A auditoria concluiu que a média do grau de conformidade dos repositórios analisados foi de 2.0 de um máximo de 5.0, tendo o ponto mais fraco sido a Sustentabilidade Financeira, com uma pontuação média de 1.39, e o ponto mais forte sido a Estrutura Governativa e Viabilidade Organizacional, com uma pontuação de 2.33 [2].

Um aspeto que distingue esta certificação da DSA, do nestor-SEAL e do CoreTrustSeal é que a única instituição que a fornece, a PTAB³⁰, não obriga as instituições responsáveis pelo repositório certificado a publicar os relatórios de certificação. Consequentemente, até 2018 não havia um relatório sobre o processo de certificação 16363:2012 de nenhum repositório disponível para análise [12].

²⁵Ligação: www.datasealofapproval.org

²⁶Ligação: <https://easy.dans.knaw.nl/ui/datasets/id/easy-dataset:116038>

²⁷Ligação: https://www.langzeitarchivierung.de/Webs/nestor/EN/Services/nestor_Siegel/nestor_siegel_node.html

²⁸Serviço de Alojamento de Repositórios Institucionais

²⁹Repositório Científico de Acesso Aberto de Portugal

³⁰Primary Trustworthy Digital Repository Authorisation Body

5.4 European Framework for Audit and Certification of Digital Repositories

A EFACDR (European Framework for Audit and Certification of Digital Repositories) surgiu após a assinatura de um memorando³¹ por parte das entidades responsáveis pelo DSA, o CCSDS³², e o standard ISO 16364:2012 e também o grupo de trabalho responsável pela DIN³³ 31644.

Esta framework propõe 3 passos de certificação: a certificação básica de acordo com o DSA, uma “Extended Certification” que inclui uma auditoria feita pela própria instituição e avaliada externamente baseada na ISO 16363 ou na DIN 31644. Por fim, o terceiro passo é certificação formal com base nas mesmas normas, mas levada a cabo através de auditoria e certificação externa [12].

5.5 CoreTrustSeal

O CoreTrustSeal, lançado em 2017, resulta do trabalho do DSA e do Repository Audit and Certification DSA–WDS Partnership Working Group da RDA. Este processo de certificação é citado nas Recomendações 9 e 13 do relatório da Comissão Europeia para a implementação dos princípios FAIR [6] como um método de certificação certificado pela comunidade que deve ser usado como base para a avaliação e certificação de serviços FAIR. Tendo em conta estas referências, e caso não existam requisitos dos stakeholders que obriguem à procura de uma certificação ISO, esta deve ser a certificação que a maioria dos repositórios institucionais deve procurar obter, pelo seu equilíbrio entre prestígio e esforço na sua obtenção. À data de escrita deste artigo existem 96 repositórios certificados de acordo com o CoreTrustSeal³⁴, entre os quais se destaca o Portulan CLARIN a nível nacional³⁵.

6 DISCUSSÃO E CONCLUSÕES

Os princípios FAIR para a gestão de dados de investigação, e mais recentemente os princípios TRUST, assumem-se como as linhas orientadoras para o desenho de fluxos de trabalho de gestão de dados, nos quais o repositório se assume como a pedra fundamental. Enquanto que os primeiros se focam nas boas práticas necessárias para a reutilização de conjuntos de dados, os segundos orientam a construção e manutenção de repositórios de dados confiáveis.

Os repositórios têm vindo a assumir um papel cada vez mais importante na divulgação de ciência. Enquanto arquivos de dados de investigação, o seu valor é cada vez mais medido pela sua capacidade de fomentar a reutilização dos recursos neles depositados. Desta forma, o simples armazenamento e descrição dos documentos neles depositados não chega. Os repositórios devem permitir a publicação de dados e metadados de forma interoperável, de forma a suportar a descoberta desses recursos, tanto por seres humanos como por máquinas, de forma automática.

Para além da descoberta e interrogação automática dos seus conteúdos por parte de sistemas externos, os repositórios devem oferecer uma interface de interação programática—ou API—completa. Só desta forma poderão suportar a execução de Planos de Gestão

de Dados Acionáveis por Máquinas, ou maDMP. Estes modelos vão além dos documentos convencionais, ao servirem de especificação legível, auditável e executável por máquinas das práticas a seguir durante a gestão de dados. Por exemplo, um maDMP pode especificar que determinados conjuntos de dados devem estar disponíveis, devem incluir um DOI nos seus identificadores e os seus metadados devem obedecer a determinada norma: o repositório onde esses dados estiverem depositados tem que ser capaz de responder com essa informação—via API—quando o sistema externo encarregado de assegurar o cumprimento do maDMP o interrogar sobre a presença desses elementos.

A capacidade de versionamento consequente auditoria transparente às alterações feitas a um conjunto de dados é também uma funcionalidade essencial nos repositórios atuais. Ambas são pré-requisito para a rastreabilidade desses dados, permitindo não só a sua evolução contínua—mesmo após a publicação como anexo a um artigo científico. Ao mesmo tempo, este registo permite atribuir aos autores dessas modificações o devido crédito.

A atribuição de crédito pelo trabalho de produção e descrição correta de conjuntos de dados é uma das questões mais relevantes para a motivação dos investigadores em todo este processo. As plataformas de repositório desempenham aqui também um papel essencial, pois são elas que devem manter a informação de quem criou ou modificou cada conjunto de dados e seus metadados que permitirá o cálculo de uma métrica de mérito. A inclusão dessa métrica de publicação nos critérios de avaliação institucional dos investigadores torna-se assim tecnicamente possível.

A oferta de capacidades de computação sobre os dados depositados num repositório requer a ligação entre os dados e plataformas de computação como por exemplo os Jupyter Research Notebooks. Desta forma, o repositório guardará não só os dados mas também os processos de análise desses dados que sustenta às conclusões publicadas. Torna-se também mais fácil para terceiros re-executar esses passos, pois reduz-se o número de dependências a instalar para recuperar o seu contexto de execução.

Por último, e em jeito de conclusão, a tecnologia desempenha um papel essencial no suporte à reprodutibilidade dos resultados apresentados nas publicações. Contudo, essa mesma tecnologia não pode nunca substituir o papel dos curadores de dados. O investimento na implementação de um repositório deve assim ser alicerçado num igual investimento na formação dos responsáveis pela gestão dos dados de investigação. São esses peritos que levam a cabo o suporte aos investigadores e que conseguem dessa forma criar um clima de confiança no processo de gestão de dados e dessa forma salientar a proposta de valor da plataforma de repositório que o sustenta.

AGRADECIMENTOS

O autor agradece o convite para participar no grupo de trabalho para a Estratégia para os Dados Abertos, e convida todos os leitores deste artigo a submeter revisões e correções no repositório GitHub deste artigo³⁶.

³¹Ligação: <http://www.trusteddigitalrepository.eu/Memorandum%20of%20Understanding.html>

³²Consultative Committee for Space Data Systems

³³Deutsches Institut für Normung

³⁴Ligação: <https://www.coretrustseal.org/why-certification/certified-repositories/>

³⁵Ligação: <https://portulanclarin.net/ecosystem/#certification>

³⁶Ligação: <https://github.com/silvae86/repositories-paper-eda>

REFERÊNCIAS

- [1] Ricardo Carvalho Amorim, João Aguiar Castro, João Rocha da Silva, and Cristina Ribeiro. 2017. A comparison of research data management platforms: architecture, flexible metadata and interoperability. *Universal Access in the Information Society* 16, 4 (2017), 851–862. <https://doi.org/10.1007/s10209-016-0475-y>
- [2] José Carvalho, Miguel Ferreira, Eloy Rodrigues, Pedro Principe, Luis Faria, Hélder Silva, and João Moreira. 2014. Auditoria ISO 16363 a repositórios institucionais. *Atas da 5ª Conferência Luso-Brasileira sobre Acesso Aberto 2* (2014). <http://hdl.handle.net/1822/30499>
- [3] Paolo Ciccarese, Stian Soiland-Reyes, Khalid Belhajjame, Alasdair JG Gray, Carole Goble, and Tim Clark. 2013. PAV ontology: provenance, authoring and versioning. *Journal of biomedical semantics* 4, 1 (2013), 37.
- [4] Thomas R Gruber. 1995. Toward principles for the design of ontologies used for knowledge sharing? *International journal of human-computer studies* 43, 5-6 (1995), 907–928.
- [5] Megan L. Head, Luke Holman, Rob Lanfear, Andrew T. Kahn, and Michael D. Jennions. 2015. The Extent and Consequences of P-Hacking in Science. *PLOS Biology* 13, 3 (03 2015), 1–15. <https://doi.org/10.1371/journal.pbio.1002106>
- [6] S Hodson, S Collins, F Genova, N Harrower, S Jones, L Laaksonen, D Mietchen, R Petrauskaitė, and P Wittenburg. 2018. Turning FAIR into reality: Final report and action plan from the European Commission expert group on FAIR data. *European Union: Brussels, Belgium* (2018).
- [7] Indian Ministry of Culture. 2017. *National Cultural Audiovisual Archives*. Technical Report May. Indian Ministry of Culture. <http://ncaa.gov.in/repository/>
- [8] Sarah Jones. 2011. How to Develop a Data Management and Sharing Plan. <https://www.dcc.ac.uk/guidance/how-guides/develop-data-plan>. Accessed: May 2020.
- [9] Shadrack Katuu. 2013. The Utility of Maturity Models—The ECM Maturity Model within a South African context. *Capability assessment and improvement workshop (CAIW) at IPRES* July (2013), 6.
- [10] Timothy Lebo, Satya Sahoo, Deborah McGuinness, Khalid Belhajjame, James Cheney, David Corsar, Daniel Garijo, Stian Soiland-Reyes, Stephan Zednik, and Jun Zhao. 2013. Prov-o: The prov ontology. *W3C recommendation* 30 (2013).
- [11] Dawei Lin, Jonathan Crabtree, Ingrid Dillo, Robert R. Downs, Rorie Edmunds, David Giarretta, Marisa De Giusti, Hervé L'Hours, Wim Hugo, Reyna Jenkyns, Varsha Khodiyar, Maryann E. Martone, Mustapha Mokrane, Vivek Navale, Jonathan Petters, Barbara Sierman, Dina V. Sokolova, Martina Stockhause, and John Westbrook. 2020. The TRUST Principles for digital repositories. *Scientific Data* 7, 1 (2020), 144. <https://doi.org/10.1038/s41597-020-0486-7>
- [12] Michelle Lindlar and Franziska Schwab. 2019. 203.5 All that work ... for what? Return on investment for trustworthy archive certification processes – a case study. <https://doi.org/10.17605/OSF.IO/8A3SC>
- [13] B. M. Marques, J. R. Da Silva, and T. Devezas. 2019. Visualization in Reproducible Science. In *2019 14th Iberian Conference on Information Systems and Technologies (CISTI)*. 1–4.
- [14] William K. Michener. 2015. Ten Simple Rules for Creating a Good Data Management Plan. *PLOS Computational Biology* 11, 10 (10 2015), 1–9. <https://doi.org/10.1371/journal.pcbi.1004525>
- [15] Tomasz Miksa, Stephanie Simms, Daniel Mietchen, and Sarah Jones. 2018. Ten simple rules for machine-actionable data management plans (preprint). <https://doi.org/10.5281/zenodo.1172673>
- [16] National Science Foundation. 2011. Grants.Gov Application Guide A Guide for Preparation and Submission of NSF Applications via Grants.gov. (2011).
- [17] Bruno Tavares, Filipe Figueiredo Correia, and André Restivo. 2020. Trusted Data Transformation with Blockchain Technology in Open Data. In *Distributed Computing and Artificial Intelligence, 16th International Conference, Special Sessions*, Enrique Herrera-Viedma, Zita Vale, Peter Nielsen, Angel Martin Del Rey, and Roberto Casado Vara (Eds.). Springer International Publishing, Cham, 213–216.
- [18] Ruben Verborgh, Miel [Vander Sande], Olaf Hartig, Joachim [Van Herwegen], Laurens [De Vocht], Ben [De Meester], Gerald Haesendonck, and Pieter Colpaert. 2016. Triple Pattern Fragments: A low-cost knowledge graph interface for the Web. *Journal of Web Semantics* 37-38 (2016), 184 – 206. <https://doi.org/10.1016/j.websem.2016.03.003>
- [19] A Whyte and S. (Eds) Allard. 2014. How to Discover Research Data Management Service Requirements. <https://www.dcc.ac.uk/guidance/how-guides/how-discover-requirements>. Accessed: May 2020.
- [20] Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific data* 3 (2016).

SIGLAS

API Application Programming Interface. 4

CSV Comma-separated Values. 3

DCC Digital Curation Centre. 4

DMP Data Management Plan. 4

DOI Digital Object Identifier. 2

EOSC European Open Science Cloud. 5

FAIR Findable, Accessible, Interoperable, Reusable. 1–3, 9

FCCN Fundação para a Computação Científica Nacional. 2

FCT Fundação para a Ciência e Tecnologia. 2

GEI Guia Estratégico de Implementação. 5

HTML HyperText Markup Language. 2

HTTP HyperText Transfer Protocol. 2

IoT Internet of Things. 5

LOD Linked Open Data. 3–5

maDMP machine-actionable Data Management Plan. 1, 4, 9

ORCID Open Researcher and Contributor ID. 3

PAV Provenance And Versioning. 3

PID Persistent Identifier. 1

RDA Research Data Alliance. 6

RDF Resource Description Framework. 2

RPRD IG Repository Platforms for Research Data Interest Group. 6

SaaS Software as a Service. 7

SPARQL SPARQL Protocol and RDF Query Language. 4, 5

TSV Tab-separated values. 3

W3C World Wide Web Consortium. 3

XML eXtensible Markup Language. 2, 3