

# A Framework for Creation of Linked Data Mashups: A Case Study on Healthcare

Gabriel Lopes  
Instituto Federal do Ceará -  
IFCE  
Fortaleza, CE, Brasil  
gabriellopes9102@gmail.com

Vânia Vidal  
Universidade Federal do  
Ceará - UFC  
Fortaleza, CE, Brasil  
vvidal@lia.ufc.br

Mauro Oliveira  
Instituto Federal do Ceará -  
IFCE  
Aracati, CE, Brasil  
amaurooliveira@gmail.com

## ABSTRACT

Linked Data promotes the publication of structured data on the Web, easing the development of an homogenized-view over heterogeneous sources, called *Linked Data Mashup view* (LDM view). But the development of this homogenized-view still is a challenging task. This article proposes a *framework* Ontology-based that aims to ease the process of creation of LDM views. This *framework* allow users without specific knowledge to create their own applications, based on their needs. We also present a case study in which we use our approach to develop an integrated view over two heterogeneous sources from Brazilian Public Health System.

## Keywords

Linked Data Mashup; Data Integration ; Semantic Integration; Semantic Mediator

## 1. INTRODUÇÃO

A iniciativa *Linked Data* [2] promove a publicação de bases de dados anteriormente isoladas como fontes RDF interligadas. O *Linked Data* trouxe novas oportunidades para o desenvolvimento de aplicações semânticas. Essas aplicações consomem os dados de um *Linked Data Mashup* (LDM), uma aplicação *web* que promove a integração de bases de dados através da combinação, agregação e transformação de fontes possivelmente heterogêneas [5]. Existem exemplos de aplicações *Linked Data Mashup* em diversos domínios, como na Saúde [1, 11] e na Música [6]. O DrugBank [11], por exemplo, é um *mashup* que integra diversas fontes de dados abertos, reunindo informações sobre mais de 5000 medicamentos, e pode ser utilizado pela comunidade para o desenvolvimento de diversas aplicações.

Entretanto, para o desenvolvimento de aplicações LDM, é necessário a obtenção de uma visão homogênea dos dados provenientes de diferentes fontes, muitas vezes com formatos e/ou *schemas* heterogêneos. De acordo com [9], o desenvolvimento de uma visão integrada, chamada de *visão Linked Data Mashup*, não é uma tarefa simples, pois requer

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WebMedia '16, November 08-11, 2016, Teresina, PI, Brazil

© 2016 ACM. ISBN 978-1-4503-4512-5/16/11... \$15.00

DOI: <http://dx.doi.org/10.1145/2976796.2988185>

a *integração semântica* das fontes de dados, i.e, as diferenças semânticas entre as fontes precisam ser conciliadas. Os principais desafios para a integração semântica no contexto de *Linked Data* estão relacionados aos seguintes fatores: (i) heterogeneidade das fontes de dados e dos vocabulários usados; (ii) qualidade dos dados que podem estar fragmentados, incompletos, incorretos ou inconsistentes; (iii) resolução de conflito de URIs, visto que diferentes URIs podem se referir ao mesmo objeto. Nesse processo, são necessários conhecimentos específicos que dificultam para um usuário geral, um gestor da área da Saúde por exemplo, criar seu próprio *mashup*.

Esse artigo propõe um *Framework* para facilitar a criação de *visões Linked Data Mashup*. Um dos objetivos da proposta, é permitir que usuários finais, com poucos conhecimentos em Computação, desenvolvam suas próprias visões integradas sobre informações distribuídas. Na nossa abordagem, uma visão LDM é especificada baseada na metodologia proposta em [9]. Essa especificação indica que as fontes de dados utilizadas para o *mashup* foram integradas semanticamente, total ou parcialmente. A partir dessa especificação, usuários finais podem construir seus próprios *mashups*, chamados de *visões de aplicação mashup*, por meio de uma *interface web*, de maneira fácil, intuitiva e sem a necessidade de conhecimentos específicos. Um dos diferenciais da nossa abordagem, é que a integração semântica é realizada uma única vez, e a partir de então, a especificação gerada será reutilizada para a criação de aplicações posteriores. Para isso, fazemos uma reescrita da especificação original, utilizando os parâmetros passados pelo usuário por meio da *interface web* gráfica. Além disso, apresentamos um estudo de caso em que o *framework* proposto é utilizado para especificar uma *visão de mashup* chamada de *Datasus.hub*, que integra semanticamente duas bases heterogêneas do Sistema Público de Saúde brasileiro. Nesse estudo de caso, também demonstramos como um usuário final pode criar seu próprio *mashup*.

O restante do artigo está organizado da seguinte forma. A Seção 2 apresenta o *framework* proposto. Na Seção 3, nossa abordagem é utilizada para construir uma visão integrada sobre duas bases da Saúde. Finalmente, a Seção 4 apresenta nossas expectativas e trabalhos futuros.

## 2. ESPECIFICAÇÃO DO FRAMEWORK

### 2.1 Visão geral

Nessa Seção, apresentamos o *framework* com 4 camadas baseado em Ontologias (Figura 1), que tem como objetivo

permitir que usuários finais obtenham uma visão integrada de informações distribuídas. Na nossa abordagem, uma *integração semântica* é realizada baseada nos conceitos de [9], brevemente discutido na Seção 1, que apresenta cinco passos para a especificação e materialização de uma visão LDM. O resultado dessa integração semântica é a especificação de uma *Visão Integrada M*, que será utilizada para a criação de *visões de aplicação mashups*, ou simplesmente *visões de mashup*. Cada visão de *mashup* é criada por um usuário final por intermédio de uma *interface web* gráfica. As quatro camadas do *framework* são descritas a seguir.

### 2.1.1 Framework 4 camadas

Na Camada de Integração Semântica, a *Ontologia de Domínio*  $O_M$  representa uma visão homogênea sobre as fontes de dados  $S_1, \dots, S_n$ , que deseja-se unificar. Na Camada de Dados, cada fonte  $S_i$  é descrita por uma ontologia  $O_{S_i}$  e exporta uma ou mais visões, chamadas de *Visões Exportadas*. Cada visão exportada  $E_i$  é composta por uma *Ontologia Exportada*  $O_{E_i}$ , cujo vocabulário é um subconjunto de  $O_M$ , e um conjunto de mapeamentos  $M_{E_i}$ , que mapeia os conceitos de  $O_{S_i}$  em  $O_{E_i}$ . Também é na Camada de Visões Exportadas que são definidas as regras para descobertas de *Links Semânticos*,  $EL_1, \dots, EL_m$ , que apontam a similaridade entre dois objetos do mundo real em duas bases distintas.

Usuários finais utilizam os parâmetros  $(O_{V_i}, F_i)$  para construir visões de *mashup* baseadas na especificação da *Visão Integrada* previamente criada. A *Ontologia de Aplicação*  $O_{V_i}$ , cujo vocabulário deve ser um subconjunto do vocabulário de  $O_M$ , representa os conceitos de interesse do usuário; enquanto  $F_i$  representa um conjunto de filtros (e.g. cidade, ano etc.) que serão aplicados sobre os dados. Os usuários finais interagem com o *framework* por meio de uma *interface web* gráfica, que permite, de forma intuitiva, a construção de uma visão de *mashup*. Como será discutido na subseção 2.3, para a materialização da visão  $V_i$ , é necessário combinar os parâmetros  $(O_{V_i}, F_i)$  com a especificação da visão integrada  $M$ , gerando uma especificação de  $V$ . Essa especificação  $V$  denota uma integração semântica parametrizada sobre as bases  $S_1, \dots, S_n$ . A especificação de  $V$  será então utilizada na materialização dos dados, realizada automaticamente com o auxílio de *frameworks* específicos.

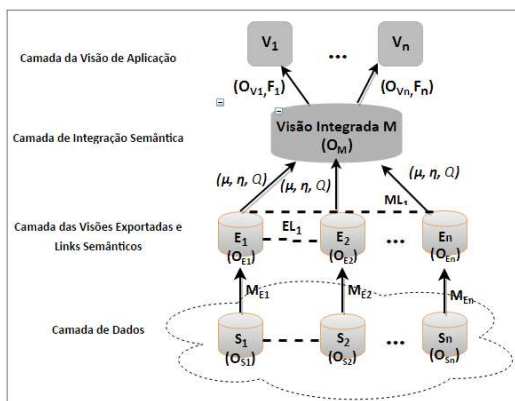


Figura 1: *Framework 4 Camadas*

## 2.2 Especificação da Visão Integrada

Seja  $M$  uma visão em um formato homogêneo sobre as fontes de dados que deseja-se unificar. A especificação de  $M$  é uma 6-tupla:

$$\bullet \lambda_m = \{M, O_M, E_M, EL_M, \mu_m, Q\}.$$

Onde  $M$  é o nome da visão integrada;  $O_M$  é a ontologia de domínio;  $E_M$  representa o conjunto das visões exportadas;  $EL_M$  define o conjunto de regras para descoberta dos *links* semânticos;  $\mu_m$  determina as regras de fusão e, finalmente,  $Q$  determina os critérios para avaliação de qualidade das fontes de dados. O processo de especificação da visão  $M$  consiste em 5 etapas:

1. **Seleção das fontes de dados relevantes para a aplicação.** As fontes de dados  $S_1, \dots, S_n$  são escolhidas de acordo com a relevância para a aplicação. Por exemplo, se a aplicação requer dados geográficos, GeoNames<sup>1</sup> pode ser uma escolha. Além disso, deve haver uma ontologia  $O_{S_i}$  descrevendo cada fonte  $S_i$ .
2. **Modelagem da Ontologia de Domínio.** Denotada por  $O_M$ , essa ontologia representa todos os conceitos do *framework* que podem ser utilizados para a criação das visões de aplicação *mashup*.
3. **Especificação das Visões Exportadas.** A especificação de uma  $E_i \in E_M$  consiste em: (i) modelar a ontologia exportada  $O_{E_i}$ , cujo vocabulário deve ser um subconjunto do vocabulário de  $O_M$  e (ii) especificar o conjunto de mapeamentos  $M_{E_i}$  que mapeia os conceitos de  $O_{S_i}$  para  $O_{E_i}$ .
4. **Especificação dos conjuntos de Links Semânticos.** Cada conjunto  $EL_i \in EL_M$  denota regras para definir que  $x$  objetos do mundo real em visões exportadas distintas,  $E_u$  e  $E_v$ , representam uma mesma entidade.
5. **Especificação as regras de fusão e de avaliação de qualidade.** As regras de fusão  $\mu_m$  explicitam como duas representações distintas de um mesmo objeto do mundo real serão combinadas em uma única representação. As regras de avaliação  $Q$  são utilizadas para quantificar a qualidade das fontes de dados.

## 2.3 Construção de uma Visão de Aplicação

Para a construção de uma Visão de Aplicação *Mashup*, denotada por  $V$  na Figura 1, são necessárias 3 etapas: (i) geração da especificação  $V$  sobre  $M$ ; (ii) geração da especificação  $V$  sobre as fontes de dados e (iii) materialização de  $V$ . Essas etapas são descritas a seguir.

### 2.3.1 Geração da Especificação de $V$ sobre $M$

A especificação de  $V$  sobre  $M$  é uma tupla  $V = (O_V, F_V)$ , onde  $O_V$  é a ontologia de aplicação, que deve ser um subconjunto da ontologia de domínio  $O_M$  e representa os conceitos de interesse do usuário; e  $F_V$  é um conjunto de filtros definidos sobre os conceitos de  $O_V$ . A especificação de  $V$  pode ser realizada através de uma *interface web* gráfica, onde, por meio de uma seleção intuitiva, o usuário define a ontologia de aplicação e os filtros.

<sup>1</sup><http://www.geonames.org/>



## 3.2 Construção de uma visão de mashup com Datasus\_hub

Nessa subseção, utilizamos a especificação de *Datasus\_hub* e as etapas definidas na Seção 2.3 para criar uma *visão de aplicação mashup*. Em nosso exemplo fictício, um gestor da saúde no Brasil quer alertar a população de seu município sobre os perigos dos maus-hábitos durante a gravidez. Para isso, o gestor quer correlacionar o uso de drogas, do tabaco e de álcool durante a gestação com a malformação em recém-nascidos. No Brasil, essas informações estão distribuídas em fontes de dados com formatos heterogêneos. Assim, utilizamos nosso *framework* para criar uma *visão de aplicação mashup V*, que atenda às necessidades do gestor, baseada na visão integrada *Datasus\_hub*. As etapas para construção da visão são apresentados a seguir.

### 3.2.1 Seleção dos Filtrros e da Ontologia de Aplicação

Nessa etapa, utilizamos a *interface* gráfica para selecionar os conceitos (classes) da ontologia de domínio (Fig. 2) importantes para nossa aplicação. Dessa seleção, extraímos a ontologia de aplicação  $O_V$ . As classes selecionadas foram: *foaf:Person*, *gissa:Mae*, *gissa:Gestacao*, *gissa:UsuarioDrogas*, *gissa:Fumante*, *gissa:Municipio*, *gissa:AnomaliaCongenita*, *gissa:Nascimento*, *gissa:UsuarioAlcool*. No nosso exemplo, o gestor quer construir uma aplicação sobre as informações de seu município. Portanto, nosso  $F_V = \{gissa:Municipio\}$ . As próximas etapas acontecem de forma transparente para o usuário.

### 3.2.2 Geração das Visões Exportadas

Nosso *framework* combina a especificação de *Datasus\_hub* com os parâmetros  $O_V$  e  $F_V$  para gerar a especificação de  $V$ . O conjunto das ontologias exportadas de  $V$ ,  $E_V$ , é definido por:  $E_V = \{(O_V \cap O_{E_{sinasc}}) \cup (O_V \cap O_{E_{sus}})\}$ , denotadas por  $E_{V_s}$  e  $E_{V_e}$  respectivamente. As regras de mapeamento e de descoberta dos *links semânticos* são importadas da especificação de *Datasus\_hub* conforme previamente definido.

### 3.2.3 Materialização da Aplicação Mashup

Nessa etapa, o *framework* retorna uma visão homogênea dos dados anteriormente isolados. Como definido em 2.3.3, a materialização ocorre em 3 etapas. Na primeira etapa, utilizamos o SILK [3] para descoberta dos *links semânticos* representados por *owl:sameAs*, resultando em 326 *links* entre  $E_{V_s}$  e  $E_{V_e}$ . Para a materialização das visões exportadas, utilizamos uma versão modificada D2R-Server [4] para processar os mapeamentos R2RML [10]. Na materialização da visão de aplicação, utilizamos o SIEVE [7] para definir as regras de qualidade e de fusão dos dados. Na avaliação de qualidade dos dados, a base  $S_{sinasc}$  obteve uma pontuação de 643, enquanto  $S_{sus}$ , 248. Desta forma, em caso de fusão, utilizamos as informações do  $S_{sus}$ .

## 4. CONCLUSÃO E TRABALHOS FUTUROS

Neste artigo, apresentamos um *framework* baseado em Ontologias para auxiliar que usuários gerais, com poucos conhecimentos específicos em Computação, possam desenvolver visões integradas sobre dados em fontes heterogêneas. Dentre os benefícios que esperamos com a nossa abordagem, podemos citar os seguintes. (1) Fácil integração de novas fontes de dados. (2) Auxiliar gestores na tomada de deci-

são. (3) A partir da integração semântica realizada pelo *framework*, quaisquer bases que utilizem os mesmos *schemas*, ou conceitos parecidos, poderão ser integradas com o mínimo de esforço. (4) Modificação mínima em casos de mudanças nos *schemas* das fontes de dados. O *framework* apresentado está na fase de conceptualização, portanto, um dos nossos trabalhos futuros é o desenvolvimento de uma ferramenta *web* que vai implementar nossa metodologia. Além disso, futuramente esperamos avaliar nossa abordagem utilizando-a para auxiliar tomadores de decisão na Saúde Pública do Brasil.

## 5. REFERENCES

- [1] F. Belleau, M.-A. Nolin, N. Tourigny, P. Rigault, and J. Morissette. Bio2rdf: Towards a mashup to build bioinformatics knowledge systems. *J. of Biomedical Informatics*, 41(5):706–716, Oct. 2008.
- [2] C. Bizer, T. Heath, and T. Berners-Lee. Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*, 5(3):1–22, 2009.
- [3] C. Bizer, J. Volz, G. Kobilarov, and M. Gaedke. Silk - a link discovery framework for the web of data. In *18th International World Wide Web Conference*, April 2009.
- [4] D2RQ. *D2RQ - Accessing Relational databases as Virtual RDF Graphs*, mar 2012. available at <http://d2rq.org/>.
- [5] H. H. Hoang, T. N. Cung, D. K. Truong, D. Hwang, and J. J. Jung. Semantic information integration with linked data mashups approaches. *IJDSN*, 2014, 2014.
- [6] G. Kobilarov, T. Scott, Y. Raimond, S. Oliver, C. Sizemore, M. Smethurst, C. Bizer, and R. Lee. *Media Meets Semantic Web - How the BBC Uses DBpedia and Linked Data to Make Connections*, pages 723–737. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009.
- [7] P. N. Mendes, H. Mühleisen, and C. Bizer. Sieve: Linked Data Quality Assessment and Fusion. In *2nd International Workshop on Linked Web Data Management (LWDM 2012) at the 15th International Conference on Extending Database Technology, EDBT 2012*, page to appear, March 2012.
- [8] L. L. Pipino, Y. W. Lee, and R. Y. Wang. Data quality assessment. *Commun. ACM*, 45(4):211–218, Apr. 2002.
- [9] V. M. P. Vidal, M. A. Casanova, N. Arruda, M. Roberval, L. P. Leme, G. R. Lopes, and C. Renso. *Advanced Information Systems Engineering: 27th International Conference, CAiSE 2015, Stockholm, Sweden, June 8-12, 2015, Proceedings*, chapter Specification and Incremental Maintenance of Linked Data Mashup Views, pages 214–229. Springer International Publishing, Cham, 2015.
- [10] W3C. *R2RML RDB to RDF Mapping Language*, jun 2016. available at <https://www.w3.org/TR/r2rml/>.
- [11] D. S. Wishart, C. Knox, A. Guo, D. Cheng, S. Shrivastava, D. Tzur, B. Gautam, and M. Hassanali. Drugbank: a knowledgebase for drugs, drug actions and drug targets. *Nucleic Acids Research*, 36(Database-Issue):901–906, 2008.