

Integrating Provenance into the Web of Data

Olaf Hartig¹ and Jun Zhao²

¹ Humboldt-Universität zu Berlin
hartig@informatik.hu-berlin.de

² University of Oxford
jun.zhao@zoo.ox.ac.uk

1 Introduction

During recent years an increasing number of data providers adopted a set of best practices for publishing and connecting structured data on the Web, leading to the creation of a globally distributed dataspace, the Web of Linked Data [1]. While this dataspace holds an enormous potential, using data from the Web poses questions of information quality and trustworthiness. Provenance information about who created and published the data and how, provides a promising means for quality assessment [2]. However, to enable successful and reliable provenance-based assessment methods we conceive it as an absolute necessity that provenance-related metadata becomes an integral part of the Web of Linked Data [3]. This is only possible if the publication of this metadata adheres to the same principles that are used for the data itself. Therefore, we present a vocabulary that allows providers of Linked Data to describe the provenance of their data with RDF. These descriptions have to be created and made available, ideally as Linked Data, again. To initiate such a practice with a minimal effort for data publishers we extended several Linked Data publishing tools with corresponding metadata components. In the remainder of this poster proposal we introduce our vocabulary in Section 2 and we describe our metadata extensions in Section 3.

2 The Provenance Vocabulary

In a recent study we reveal a lack of suitable vocabularies to describe provenance of Linked Data [4]. To fill this void we develop the Provenance Vocabulary³. We designed the vocabulary very closely to our model for Web data provenance [4]. As the provenance model comprises two dimensions, i.e., data creation and data access, our vocabulary accordingly consists of three categories of terms: general terms, terms for data creation, and terms for data access (cf. Figure 1). These terms allow to describe all common cases related to Linked Data publishing such as: manual data creation, Linked Data interfaces over native RDF stores, wrappers over relational databases or over Web APIs, file-based data creations, data changes, retrieval and provenance of source data, publication responsibilities, and service operators. To allow for a wide range of applications the vocabulary does not prescribe a specific granularity by which provenance has to be described. Hence, the classes are quite general. For instance, a DataItem could

³ <http://purl.org/net/provenance/>

