

nanopub-java: A Java Library for Nanopublications

Tobias Kuhn^{1,2}

¹ Department of Humanities, Social and Political Sciences, ETH Zurich, Switzerland

² Department of Computer Science, VU University Amsterdam, Netherlands

kuhntobias@gmail.com

Abstract. The concept of nanopublications was first proposed about six years ago, but it lacked openly available implementations. The library presented here is the first one that has become an official implementation of the nanopublication community. Its core features are stable, but it also contains unofficial and experimental extensions: for publishing to a decentralized server network, for defining sets of nanopublications with indexes, for informal assertions, and for digitally signing nanopublications. Most of the features of the library can also be accessed via an online validator interface.

1 Introduction

This technical paper presents `nanopub-java`, which is a library for nanopublications. Its source code can be found here:

<https://github.com/Nanopublication/nanopub-java>

Nanopublications³ [2,8] are an approach to publish scientific data and meta-data in RDF by subdividing them into small data snippets. They are a concrete proposal to implement the visions of *semantic publishing* [9] and *linked science* [3] by allowing for the publication and sharing of formally represented scientific resources and data that are semantically interlinked and provide provenance and context information for their reliable integration and evaluation. Specifically, a nanopublication consists of three named graphs of RDF triples (plus a fourth graph to keep them together): the assertion graph contains the actual content of the nanopublication (e.g. a scientific finding); the provenance graph contains information about the provenance of the assertion (e.g. the scientific method with which the assertion was derived); and the publication information graph contains meta-data about the nanopublication itself (e.g. its creator and a timestamp).

The library presented here can be useful in a number of scenarios:

- To represent and share small chunks of scientific knowledge and metadata in RDF in a provenance-aware manner (as nanopublications)
- To make RDF content verifiable and immutable (with trusty URIs)

³ <http://nanopub.org>

- To define large or small datasets of RDF content where the data entries can be individually addressed and recombined in new datasets (with nanopublication indexes)
- To quickly publish RDF snippets in a verifiable and permanent manner (relying on an existing server network)
- To retrieve existing nanopublications from the network (5 millions and counting)
- To digitally sign RDF snippets (though this is still experimental)

Below the details of the library and its web interface are explained.

2 Implementation

The library is built upon the Sesame library [1] to validate, represent, and create RDF structures. The features of the nanopublication library are centered around a Java interface representing a nanopublication, called `Nanopub`, and the Java class `NanopubImpl` provides a reference implementation of this interface. This implementation checks the well-formedness of a nanopublication at the time of its creation based on the latest version of the nanopublication guidelines,⁴ and raises an exception in the case of a violation of these rules.

Trusty URIs [6,7] are the recommended way of how to make nanopublications verifiable and immutable, and to give them identifiers based on cryptographic hash values. The nanopublication library uses for that purpose the `trustyuri-java` library.⁵ In a nutshell, a trusty URI is a kind of URI reference that contains a cryptographic hash value that is calculated on the digital artifact it represents. This allows one to verify that a given content is really what the URI was supposed to represent by its creator, and thereby to enforce the immutability of digital artifacts such as nanopublications.

The features of the library are made available through the Java API as well as via a command line interface using the command `np`. The following features are part of the core of the library, which means that they deal with stable and agreed-upon structures as defined by the community:

- `check` / `CheckNanopub` reads a nanopublication or several of them and checks whether any of the well-formedness criteria are violated. If a trusty URI or a digital signature is found (see below), these are checked too.
- `mktrusty` / `MakeTrustyNanopub` takes a nanopublication that does not yet have a trusty URI and transforms it into one that is identified by a newly created trusty URI.
- `fix` / `FixTrustyNanopub` takes a nanopublication with a broken trusty URI and fixes it, i.e. assigns it a new trusty URI. This is useful when a nanopublication has to be changed, which invalidates the hash. Running this command creates a *new* nanopublication with a valid trusty URI. (Nanopublications

⁴ <http://nanopub.org/guidelines/working.draft/>

⁵ <https://github.com/trustyuri/trustyuri-java>

are immutable, so changing something necessarily leads to a new nanopublication.)

In addition to these core features, the library also contains a number of unofficial extensions (which may or may not become official at some point). There is code to validate informal assertions specified as AIDA sentences [4]; code for creating index nanopublications to define small or large sets of nanopublications [5]; code for publishing and retrieving nanopublications from a decentralized server network [5]; and experimental code for digitally signing nanopublications. These features are accessible via the following commands and classes:

- **mkindex** / **MakeIndex** takes a list of nanopublications and creates an index that refers to them. A nanopublication index therefore represents a (possibly large) set of nanopublications. Such indexes are themselves formatted as nanopublications, and can therefore also be published to the server network (see below).
- **publish** / **PublishNanopub** uploads a given nanopublication that has a trusty URIs to the server network. Such a nanopublication is then distributed among the servers of the network (currently five) and made available even if some of the servers should be inaccessible at a certain point in time. In this way, the nanopublication is made permanent and its publication cannot be undone.
- **get** / **GetNanopub** reliably retrieves a given nanopublication from the decentralized server network. Nanopublications are verified according to their trusty URI, and only verified nanopublications are returned by this command. For nanopublication indexes, the whole set of nanopublications that is defined by the index can be downloaded.
- **status** / **NanopubStatus** checks whether and how often a given nanopublication (identified by its trusty URI) is found on the server network.
- **server** / **GetServerInfo** returns some information about a given server in the network, such as the number of nanopublications it contains.
- **mkkeys** / **MakeKeys** creates a new key-pair to be used to sign nanopublications.
- **sign** / **SignNanopub** takes a nanopublication and signs it with a given private key.

3 Examples

Below, some of the most important commands are explained by examples based on the **np** command line tool. The same functionality is also available via the Java API. For the sake of these examples, let us assume that we have a file called **nanopubfile.trig** that starts with the following RDF prefixes:

```
@prefix xsd: <http://www.w3.org/2001/XMLSchema#>.
@prefix dc: <http://purl.org/dc/terms/>.
@prefix pav: <http://purl.org/pav/>.
@prefix prov: <http://www.w3.org/ns/prov#>.
@prefix np: <http://www.nanopub.org/nschema#>.
@prefix ex: <http://example.org/>.
@prefix : <http://example.org/np1#>.
```

The definition of the first nanopublication in this file starts with the head graph that defines the structure of the nanopublication by linking to the other graphs:

```
:Head {
  : a np:Nanopublication; np:hasAssertion :assertion;
    np:hasProvenance :provenance; np:hasPublicationInfo :pubinfo.
}
```

The actual claim of the nanopublication is stored in the assertion graph:

```
:assertion {
  ex:drugA ex:treats ex:diseaseB.
}
```

The provenance and publication info graphs provide meta-information about the assertion and the entire nanopublication, respectively:

```
:provenance {
  :assertion prov:wasDerivedFrom ex:some_publication.
}
:pubinfo {
  : pav:createdBy <http://orcid.org/0000-0002-1267-0234>.
  : dc:created "2015-08-18T15:36:22+01:00"^^xsd:dateTime.
}
```

The lines above constitute a very simple but complete nanopublication. To make this example a bit more interesting, let us assume that our file contains two more nanopublications that have different assertions but are otherwise identical:

```
@prefix : <http://example.org/np2#>.
...
:assertion {
  ex:Gene1 ex:isRelatedTo ex:diseaseB.
}
...

@prefix : <http://example.org/np3#>.
...
:assertion {
  ex:Gene2 ex:isRelatedTo ex:diseaseB.
}
...
```

To check and validate these three nanopublications, we can now use the following command:

```
$ np check nanopubfile.trig
Summary: 3 valid (not trusty);
```

These nanopublications can now be transformed into ones with trusty URIs using the following command (resulting in a new file `trusty.nanopubfile.trig`):

```
$ np mktrusty nanopubfile.trig
```

Using the same command in verbose mode with the argument `-v` shows us the newly generated trusty URIs for the three nanopublications:

```
$ np mktrusty -v nanopubfile.trig
Nanopub URI: http://example.org/np1#RAHGBOWzgQijR88g_rIwtPCmzYgy04wRMT7M91ouhojsQ
Nanopub URI: http://example.org/np2#RA4xTdhe2gPctqvAwdgTU4eRiR1aTQ1TYJcF3Sohe5Cus
Nanopub URI: http://example.org/np3#RAEjvXP0xTkeIa2mKmYT66i_PAJ-u-k0uRBd6_sMe9qG0
```

As they are tiny snippets of data, nanopublications are most useful when they are grouped and combined in small or large collections. We therefore need a simple

method to refer to collections or sets of nanopublications, which is achieved by the experimental proposal of *nanopublication indexes* [5], which are themselves nanopublications. Such indexes can be used to group nanopublications that have trusty URIs using the following command:

```
$ np mkindex trusty.nanopubfile.trig
Index URI: http://np.inn.ac/RAFa_x4h0ng_NXtof35Ie9pQVsAY69Ab3ZQMir2NP8vGc
```

The nanopublications of the new index are saved in a file called `index.trig` unless specified otherwise with the argument `-o`.

Moving to the part that involves the server network, nanopublications that have trusty URIs (which includes nanopublication indexes) can be published to the network with the following command:

```
$ np publish trusty.nanopubfile.trig
3 nanopubs published at http://np.inn.ac/
```

The publication status of a given nanopublication can be checked like this:

```
$ np status -a http://example.org/np1#RAHGB0WzgQijR88g_rIwtPCmzYgy04wRMT7M91ouhojsQ
URL: http://np.inn.ac/RAHGB0WzgQijR88g_rIwtPCmzYgy04wRMT7M91ouhojsQ
URL: http://ristretto.med.yale.edu:8080/nanopub-server/RAHGB0WzgQijR88g_rIwtPCmzYgy0...
URL: http://nanopub-server.ops.labs.vu.nl/RAHGB0WzgQijR88g_rIwtPCmzYgy04wRMT7M91ouhojsQ
URL: http://nanopubs.stanford.edu/nanopub-server/RAHGB0WzgQijR88g_rIwtPCmzYgy04wRMT7...
URL: http://nanopubs.semantic-science.org/RAHGB0WzgQijR88g_rIwtPCmzYgy04wRMT7M91ouhojsQ
Found on 5 nanopub servers.
```

A given nanopublication that is published on the server network can be retrieved via its URI:

```
$ np get http://www.tkuhn.ch/be12nanopub/RAhV9IpiUEjbentzGivp1Lbx0BVegp5sgE3BwS0S2RAYM
```

All the servers in the network are checked until the nanopublication is found and successfully verified. This command is therefore reliable even if one or several servers are down. Instead of the complete URI, it is also possible to just specify the trusty URI artifact code:

```
$ np get RAhV9IpiUEjbentzGivp1Lbx0BVegp5sgE3BwS0S2RAYM
```

To get the content of a nanopublication index (and not just the top-most index nanopublication), argument `-c` can be used:

```
$ np get -c -o content.trig RAtF0ivB9B8cb-u3K_zElgmRBxiDwfym1yVBRY6VAyWvE
```

Argument `-o` specifies again the name of the output file. The remaining commands as introduced above are equally intuitive to use. Just entering the command without any arguments will output a list of all argument options.

4 Web Interface

Many of the features described above are made available through the nanopublication validator Web interface,⁶ an instance of which can be accessed at <http://nanopub.inn.ac>. Figure 1 shows a screenshot. With this interface, the well-formedness of nanopublications can be checked as well as the adherence to a

⁶ <https://github.com/tkuhn/nanopub-validator>

Validator (and more) for Nanopublications

This validator interface for nanopublications is provided by Tobias Kuhn. It is based on the `nanopub-java` library and uses features from `trustyuri-java` for handling Trusty URIs. It is work in progress and might contain bugs. The source code can be found on [GitHub](#).



The screenshot shows a web interface for validating nanopublications. At the top, there are several tabs: "Direct Input", "Load Example", "Upload File", "from URL", "from SPARQL Endpoint", and "from Nanopub Server". Below the tabs is a text area with the instruction "Paste/edit your nanopublication below (UTF-8 encoded)". The text area contains the following RDF code:

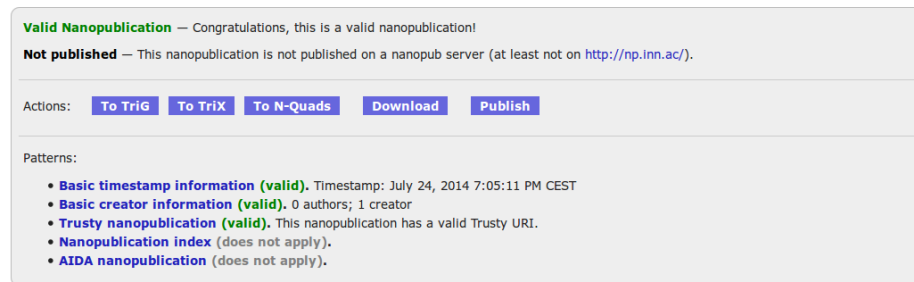
```
@prefix this: <http://example.org/nanopub-validator-example/RATAUGU_xKTH016Eoiu11Ssw0kBu1eL8_3_BoDJWH3arA#> .
@prefix sub: <http://example.org/nanopub-validator-example/RATAUGU_xKTH016Eoiu11Ssw0kBu1eL8_3_BoDJWH3arA#> .
@prefix xsd: <http://www.w3.org/2001/XMLSchema#> .
@prefix dc: <http://purl.org/dc/terms/> .
@prefix pav: <http://purl.org/pav/> .
@prefix prov: <http://www.w3.org/ns/prov#> .
@prefix np: <http://www.nanopub.org/nschema#> .
@prefix npx: <http://purl.org/nanopub/x/> .
@prefix ex: <http://example.org/> .

sub:Head {
  this: np:hasAssertion sub:assertion ;
        np:hasProvenance sub:provenance ;
        np:hasPublicationInfo sub:pubinfo ;
        a np:Nanopublication .
}
```

Below the text area is a "Check" button. In the top right corner, there is a "Format:" label with three radio buttons: "TriG" (selected), "TriX", and "N-Quads".

Nanopublication

Example loaded:



The screenshot shows the validation results for a nanopublication. It features a green header "Valid Nanopublication — Congratulations, this is a valid nanopublication!" and a red header "Not published — This nanopublication is not published on a nanopub server (at least not on <http://np.inn.ac>).". Below the headers, there are five buttons: "To TriG", "To TriX", "To N-Quads", "Download", and "Publish". Under the "Actions:" label, these buttons are visible. Below the buttons, there is a "Patterns:" section with a list of five items:

- **Basic timestamp information (valid)**. Timestamp: July 24, 2014 7:05:11 PM CEST
- **Basic creator information (valid)**. 0 authors; 1 creator
- **Trusty nanopublication (valid)**. This nanopublication has a valid Trusty URI.
- **Nanopublication index** (does not apply).
- **AIDA nanopublication** (does not apply).

Fig. 1. This is a screenshot of the nanopublication validator interface at <http://nanopub.inn.ac>.

number of patterns. They can furthermore be transformed into different RDF serializations, and published to the server network. Loading of nanopublications is possible via form input, upload, fetching from a URL, SPARQL endpoint access, and retrieval from the server network. In general, this web interface and the underlying library are supposed to support the development of best practices for the nanopublication community by providing a solid basis for discussion, by allowing for experimental features to be tested and discussed, and by facilitating the implementation of prototypes.

The current server network on which many of the unofficial features depend, can be explored via a monitor interface at <http://npmonitor.inn.ac>.

5 Conclusions

The `nanopub-java` library provides a stable implementation of the nanopublication concept, adhering to its specified guidelines. It is openly licensed under

the terms of the MIT license, is available on The Central Repository,⁷ and has so far been used in about a dozen open-source codebases.⁸

In general, we believe that this library can be a valuable resource for tools that use RDF data in the context of provenance recording, reproducibility, data publishing, data reuse, and reliable retrieval of Linked Data.

References

1. J. Broekstra, A. Kampman, and F. Van Harmelen. Sesame: A generic architecture for storing and querying rdf and rdf schema. In *The Semantic Web — ISWC 2002*, pages 54–68. Springer, 2002.
2. P. Groth, A. Gibson, and J. Velterop. The anatomy of a nano-publication. *Information Services and Use*, 30(1):51–56, 2010.
3. T. Kauppinen, A. Baglatzi, and C. Kefler. Linked science: interconnecting scientific assets. *Data Intensive Science*, 2012.
4. T. Kuhn, P. E. Barbano, M. L. Nagy, and M. Krauthammer. Broadening the scope of nanopublications. In *The Semantic Web: Semantics and Big Data — ESWC 2013*, pages 487–501. Springer, 2013.
5. T. Kuhn, C. Chichester, M. Krauthammer, and M. Dumontier. Publishing without publishers: a decentralized approach to dissemination, retrieval, and archiving of data. In *The Semantic Web — ISWC 2015*. Springer, 2015.
6. T. Kuhn and M. Dumontier. Trusty URIs: Verifiable, immutable, and permanent digital artifacts for linked data. In *The Semantic Web: Trends and Challenges — ESWC 2014*, pages 395–410. Springer, 2014.
7. T. Kuhn and M. Dumontier. Making digital artifacts on the web verifiable and reliable. *IEEE Transactions on Knowledge and Data Engineering*, 27(9), 2015.
8. B. Mons, H. van Haagen, C. Chichester, J. T. den Dunnen, G. van Ommen, E. van Mulligen, B. Singh, R. Hooft, M. Roos, J. Hammond, et al. The value of data. *Nature genetics*, 43(4):281–283, 2011.
9. D. Shotton. Semantic publishing: the coming revolution in scientific journal publishing. *Learned Publishing*, 22(2):85–94, 2009.

⁷ <https://search.maven.org/#artifactdetails|org.nanopub|nanopub|1.7|jar>

⁸ <https://github.com/Nanopublication/nanopub-java#usage-tracking>