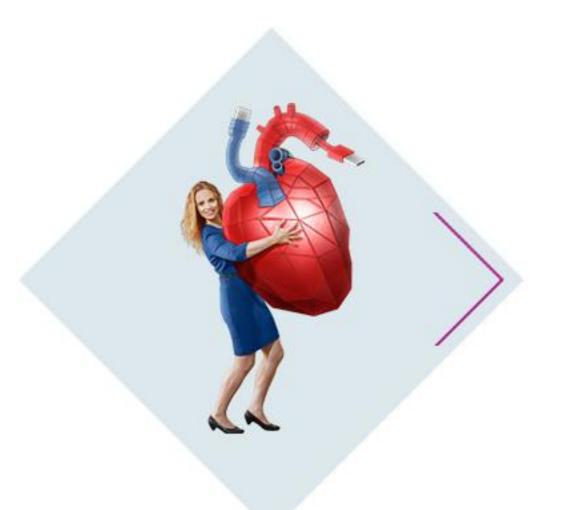


### Content

- Context: Glossary, XML Schema Definition (XSD) and external sources
- A mini-vocabulary for DmfA declarations
- Translating from DmfA XML to RDF
- SHACL rules
- Demonstration
- Summary

<u>Attention</u>: We are not going to replace DmfA and the glossary. The intention is to demonstrate the possibilities of SHACL for the validation of DmfA declarations. DmfA offers us a very complex matter with which we can demonstrate all the functionalities of SHACL.

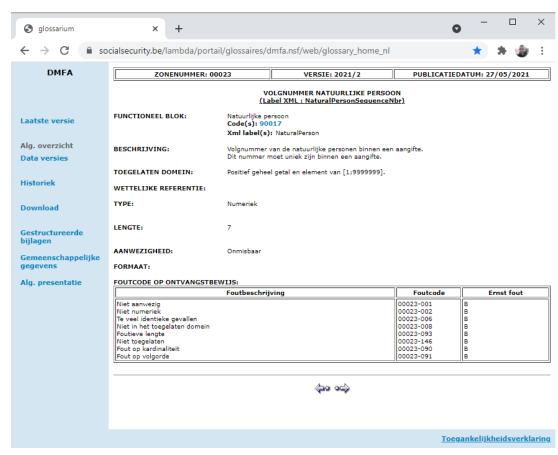




# Context

### Glossarium

- Information about entities (called blocks), relationships and attributes in HTML, PDF, XML, and TXT
- The glossary contains information on validation rules.
   Simple rules about totality (must exist) and length are structured, but many rules are described in natural language and fall outside the expressiveness of XML.
   These must be tested by an application.
- In my humble opinion: encoding bias. The model is aimed at producing an XSD.

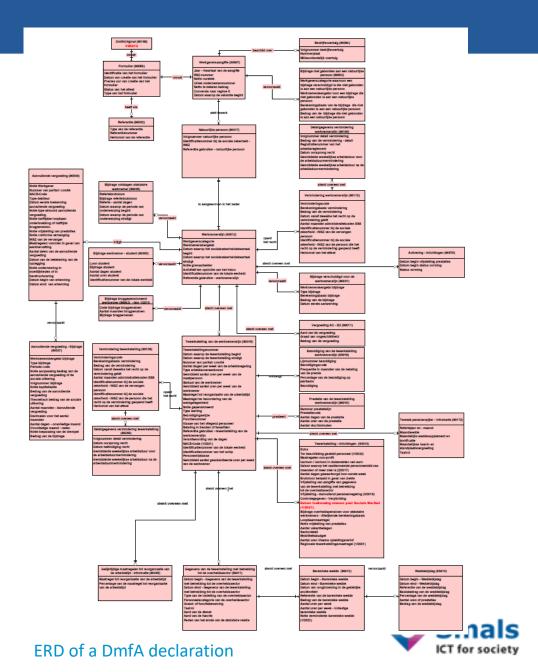


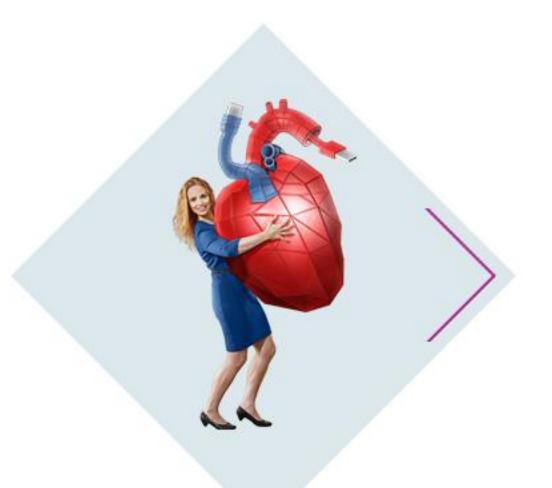
NaturalSequenceNbr is an attribute of the entity person. Sequence numbers are unique within a declaration, but this is not tested by XSD. An application must test this.



# Glossarium

- 29 entities (functional blocks)
- 28 relationships between entities (hierarchy of the blocks)
- 215 attributes (but not all is used)
- Some attributes are historical and documented, but their use is "forbidden"





# A vocabulary for DmfA declarations

## Why an ontology vocabulary? (\*)

- A direct mapping is a direct translation of data to RDF. The resulting RDF reflects the structure of the source, in this case the XML.
- Relationships in XML are implicit (see e.g. relationship between Form and Reference).

```
1 k?xml version="1.0" encoding="UTF-8"?>
 2 <!-- edited with XMLSpy v2019 cel. 3 (x64) (http://www.altova.com) by (Smals) -->
 3@ <DmfAOriginal xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:noNamespaceSchemaLocation="DmfAOriginal_20212.xsd">
       <Form>
                                                               Attribute with value of entity Form
            <Identification>DMFA</Identification> ◀
            <FormCreationDate>2021-05-28</FormCreationDate>
            <FormCreationHour>08:14:45.000/FormCreationHour>
            <AttestationStatus>0</AttestationStatus>
            <TypeForm>SU</TypeForm>
                                                               Relationship between Form and Reference: here the name of the type is also the
            <Reference>
11
                <ReferenceType>1</ReferenceType>
                                                               relationship.
                <ReferenceOrigin>1</ReferenceOrigin>
12
13
                <ReferenceNbr>14506543920212</ReferenceNbr>
14
            </Reference>
15<sup>(1)</sup>
           <EmployerDeclaration>
16
                <Quarter>20212</Quarter>
17
                <NOSSRegistrationNbr>145065439</NOSSRegistrationNbr>
```

Although a direct mapping with SHACL is possible, the SHACL rules would not be intuitive: "Do we
think in functions of entities or relationships?" Consequently, a simple vocabulary seemed appropriate
to me.



## Generating a vocabulary: entities

Glossary contains XML documentation that we can process automatically

**Using a Jupyter Notebook (see GitLab)** 

#### Creating **types** with:

- BFLabelsXMLdmfa20212FR.xml
- FBLabelsXMLdmfa20212NL.xml

Documentation in NL, FR, and EN.

```
ont:CompanyVehicle a owl:Class;

rdfs:label "Company Vehicle"@en,

"Véhicule de société"@fr,

"Bedrijfsvoertuig"@nl;

dc:identifier "90294".
```



### Generating a vocabulary: attributes

Glossary contains XML documentation that we can process automatically

#### Creating <u>attributes</u> with:

- ZonesLabelsXMLdmfa20212FR.xml
- ZonesLabelsXMLdmfa20212NL.xml

Documentation in NL, FR, and EN.

Encoding errors in XML require some preprocessing.



## Generating a vocabulary: entity-attributes

Glossary contains XML documentation that we can process automatically

# Relationships between entities and attributes:

VersCplt\_BFdmfa20212FR.xml

VersCplt FBdmfa20212NL.xml

Only the relationships needed.



## Generating a vocabulary: entity-attributes

Glossary contains XML documentation that we can process automatically

#### **Documentation of the attributes:**

- VersCpltdmfa20212FR.xml
- VersCpltdmfa20212NL.xml

Documentation in NL and FR.

```
ont:ReducedScaleSalaryNotion a owl:DataProperty;
   rdfs:label "Reduced Scale Salary Notion"@en,
        "NOTION TRAITEMENT BARÃMIQUE RÃDUIT"@fr,
        "NOTIE VERMINDERDE BAREMIEKE WEDDE"@nl ;
   rdfs:domain ont:ScaleSalary ;
    dc:identifier "01235";
    rdfs:comment "Zone qui ... réduit."@fr, "Zone ... is."@nl ;
    rdfs:isDefinedBy "1 = traitement barémique réduit"@fr,
        "1 = verminderd baremieke wedde"@nl;
    ont: function "Traitement barémique"@fr, "Baremieke wedde"@nl;
    ont:length 1 ;
    ont:presence "Obligatoire si .... "@fr, "Verplicht indien .... "@nl;
    ont:type "Alphanumérique"@fr, "Alfanumeriek"@nl ;
```



## Generating a vocabulary: relationships

We processed relationships between entities manually using the ERD.

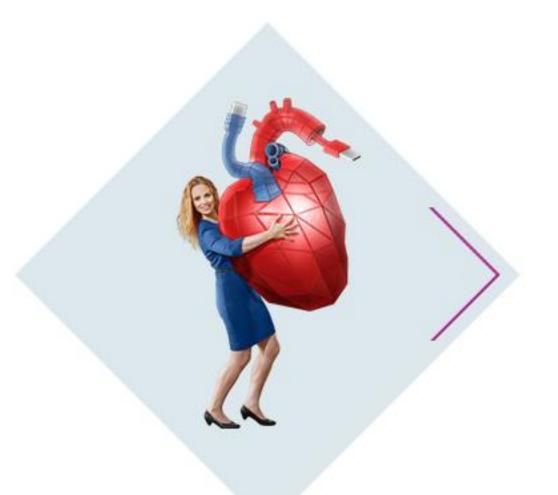
We could use the XSD, but relationships in XSD had no names. Furthermore, there are relationships with the same names. In the long run, it is better to make a distinction (at least at the level of the vocabulary).

A CSV file was created, and it was transformed:

FROM	то	LABELNL	LABELFR	LABELEN
90169	90059	omvat	reprent	contains
90059	90082	heeft als	a comme	has
90059	90007	omvat	comprent	contains
90007	90294	beschikt over	dispose de	has
			•••	•••

Note that I combine the codes of two entities for the computer label of the relationship. In this way, different "has" relationships are different.





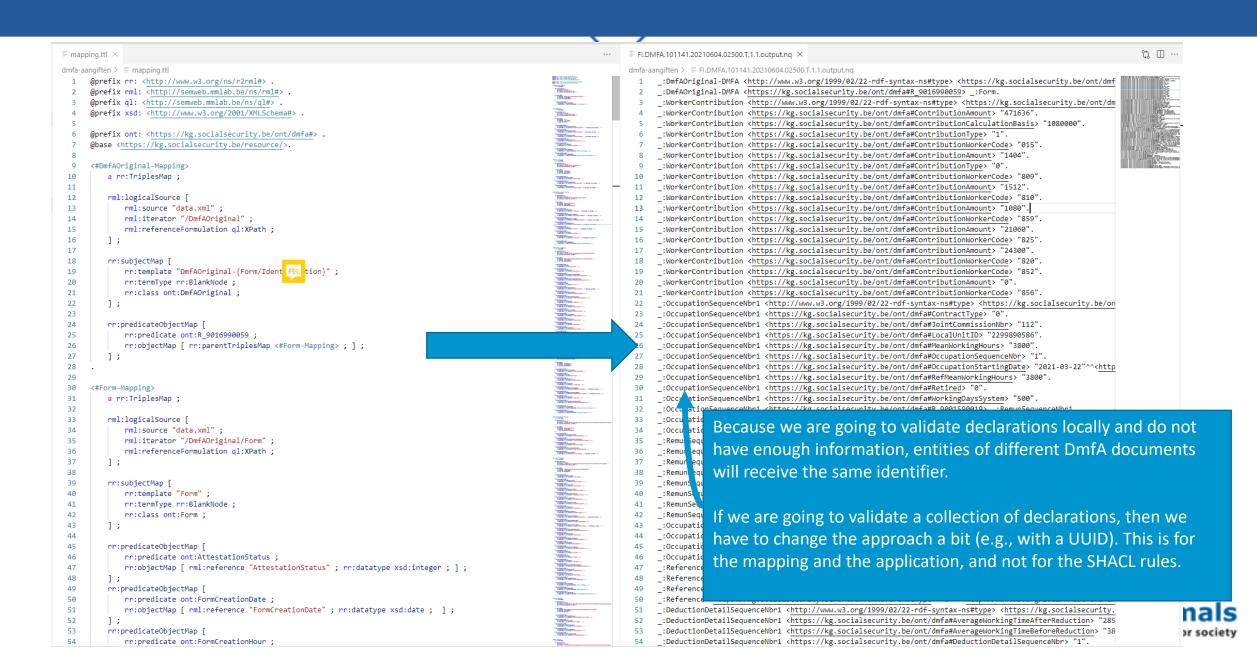
# From a DmfA XML declaration to RDF

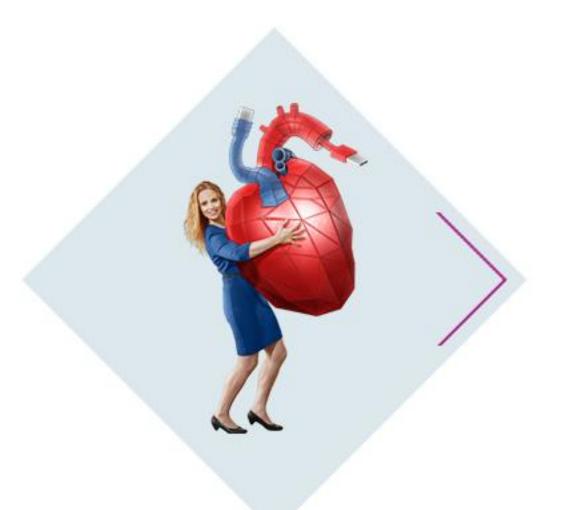
### From XML to RDF

- Based on the XML file FI.DMFA.101141.20210604.02500.T.1.1 (supplied by Smals)
- Used 11 of the 29 entities. So not everything was mapped.
- Mapping using RML (<a href="https://rml.io/">https://rml.io/</a>) the "successor" of R2RML for relational databases, CSV, JSON, XML,...
- The mapping is simple, it is almost a one-on-one mapping. Data types such as xsd:string, xsd:integer, xsd:date,... must be supplied.
- To be supplemented here and there when SHACL rules are completed.



### From XML to RDF





# **SHACL rules**

# Generating SHACL

- The generation of SHACL from the glossary is limited (type, length, etc.). There was enough for one form of a skeleton but:
  - •
  - Contradictions between glossary and rules
  - Sometimes errors in the XSD (?!)
  - Inconsistent approach to some constraints (such as length—exact length vs. maximum length)
- The skeleton therefore contains errors (beyond my abilities) and everything has to be tested and refined manually. I only gained a little time because of this... So I removed the code for the skeleton. This process can be better aligned with the generation of the XSD.



kg:DmfAOriginalCountInstancesShape: a document must contain exactly 1 entity of the type ont:DmfAOriginal, otherwise we do not have a declaration.



kg:DmfAOriginalShape: a ont:DmfAOriginal must contain at least 1 ont:Form.

```
kg:DmfAOriginalCountInstancesShape a sh:NodeShape;
    rdfs:comment "A graph should have exactly one instance of ont:DmfAOriginal";
   sh:targetNode ont:DmfAOriginal ;
    sh:property [
        sh:message "Graph does not contain exactly one instance of ont:DmfAOriginal";
        sh:path [ sh:inversePath rdf:type ] ;
        sh:maxCount 1;
        sh:minCount 1;
   ] ;
. ] ;
kg:DmfAOriginalShape a sh:NodeShape ;
   rdfs:comment "Property Shape for DmfAOriginal (code)";
   sh:targetClass ont:DmfAOriginal;
    sh:property [
        rdfs:comment "ont:DmfAOriginal must have at least one ont:Form.";
        sh:minCount 1;
        sh:path ont:R_9016990059
```

Smals ICT for society

### SHACL: voorbeelden

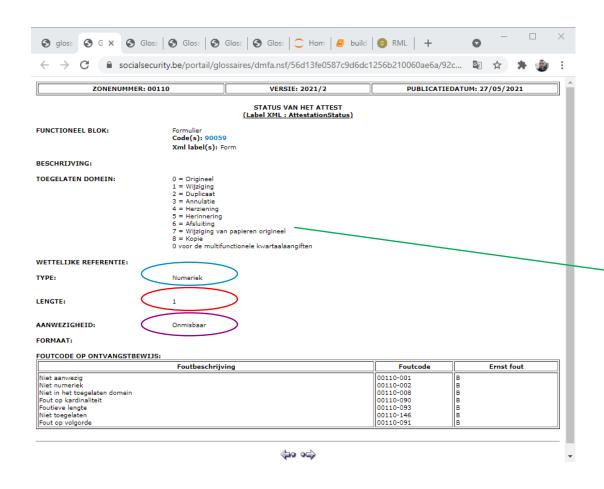
#### kg:FormShape

- Constraints on relationships with DmfAOriginal, Reference, and EmployerDeclaration
- Constraints over the attributes

We now choose a number of interesting attributes!

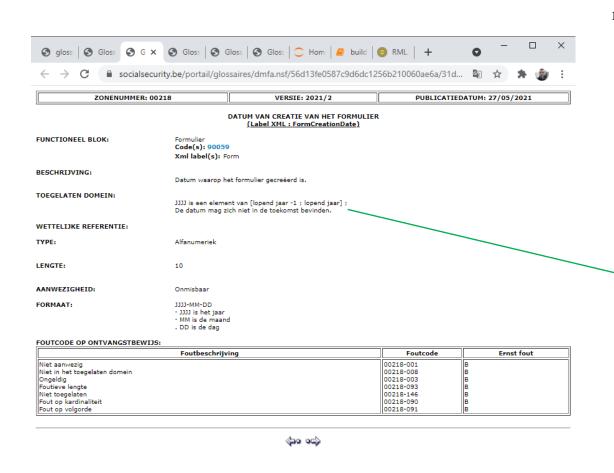
```
kg:FormShape a sh:NodeShape;
   rdfs:comment "Property Shape for Form (code)";
   sh:targetClass ont:Form ;
   sh:property [
       rdfs:comment "Each ont:Form must belong to exactly one ont:DmfAOriginal.";
       sh:maxCount 1 ; sh:minCount 1 ;
       sh:path [ sh:inversePath ont:R 9016990059 ; ] ;
   ] ;
   sh:property [
       rdfs:comment "ont:Form has at most one ont:Reference.";
       sh:maxCount 1 ;
       sh:path ont:R 9005990082;
   ] ;
   sh:property [
       rdfs:comment "Each ont:Form has exactly one ont:EmployerDeclaration.";
       sh:maxCount 1 ; sh:minCount 1 ;
       sh:path ont:R 9005990007;
   ] ;
   sh:property kg:AttestationStatusShape ;
   sh:property kg:FormCreationDateShape;
   sh:property kg:FormCreationHourShape;
   sh:property kg:IdentificationShape ;
   sh:property kg:TypeFormShape ;
```





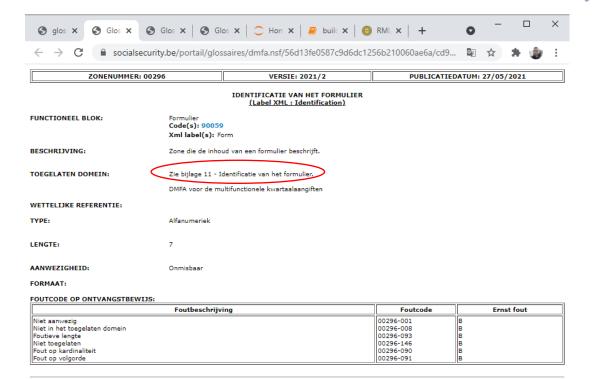
```
kg:AttestationStatusShape a sh:PropertyShape;
    rdfs:comment "Property Shape for AttestationStatus
(code)";
    sh:path ont:AttestationStatus;
    sh:datatype xsd:integer;
    sh:minCount 1;
    sh:maxCount 1;
    sh:minLength 1;
    sh:maxLength 1;
    sh:minInclusive 0;
    sh:maxInclusive 8;
```





In this preview, the sh:pattern is actually superfluous, because this is a condition of the xsd:date datatype.

```
kg:FormCreationDateShape a sh:PropertyShape;
    rdfs:comment "Property Shape for FormCreationDate (code)";
    sh:path ont:FormCreationDate ;
    sh:datatype xsd:date ;
                                                     We use the
    sh:maxLength 10;
                                              expressiveness of the
    sh:minLength 10;
                                              graph query language
    sh:pattern \sqrt{d\{4\}}-\sqrt{d\{2\}}-\sqrt{d\{2\}};
    sh:sparql
      sh:message "Year beyond [current year - 1, current year].";
      sh:prefixes <> ;
      sh:select """
       SELECT $this (ont:FormCreationDate AS ?path) ?value
        WHERE {
          $this ont:FormCreationDate ?value .
               ( YEAR ( ?value ) AS ?y1 )
               ( YEAR ( NOW () ) AS ?y2 )
          FILTER ( !( ?y2 - ?y1 = 0 || ?y2 - ?y1 = 1 ) )
```



S Bijlagen - Google Chrome socialsecurity.be/portail/glossaires/bijlagen.nsf/web/Bijlagen\_Home\_NI **BIJLAGEN** Identificatie van het formulier Laatste versie 9 Functienummers 10 Codificatie vergoedingen 12 Identificatie vh risico 13 Code aard van de dag 14 Afwijkende gebeurtenisser 15 Betrokken voorwerpen Code Omschrijving Begindatum geldigheid Einddatum geldigheid AADD501 Aanvraag aangifte 01-01-1900 01-01-9999 ACRF001 Ontvangstbewijs 01-01-1900 01-01-9999 AOAT001 Arbeidsongevallen scenario 1 - aangifte van een 01-01-1900 01-01-9999 Gewijzigde bijlagen AOAT002 Arbeidsongevallen scenario 2 - maandelijkse 01-01-1900 01-01-9999 rapportering Per status AOAT003 Arbeidsongevallen scenario 3 - mededeling van 01-01-1900 01-01-9999 AVDDTUP Wijzigende aangifte met betrekking tot een aangifte 01-01-1900 01-01-9999 Overzicht AVWDDT Aangifte van werken 01-01-9999 BEWARE Notificatie die de boekhoudkundige informatie bevat 01-01-1900 01-01-9999 in verband met de wijzigingen van de multifunctionele

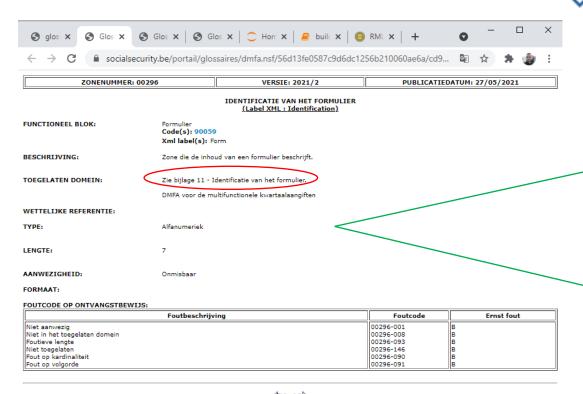
This is an example of how the different identification types in RDF can be represented:

```
kg:IdentificationTypeDMFA a ont:IdentificationType;
    rdfs:label "DMFA";
    rdfs:comment "Multifunctionele ... werkgever";
    ont:validFrom "1900-01-01"^^xsd:date;
    ont:validTo "9999-01-01"^^xsd:date;
```

With such representations in our data graph (or even knowledge graph) we can:



### SHACL: voorbeelden

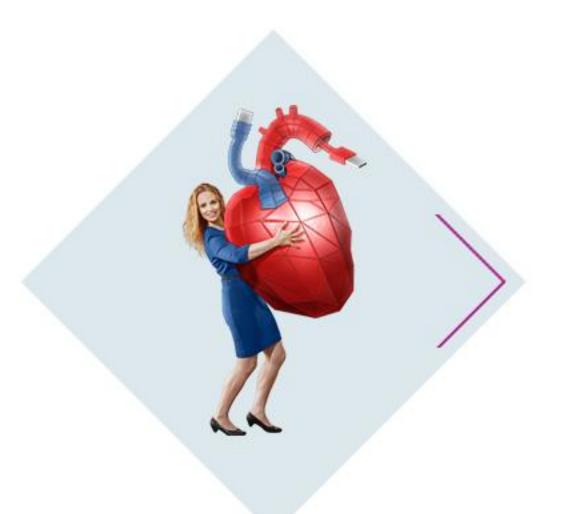


```
3 Bijlagen - Google Chrome

■ socialsecurity.be/portail/glossaires/bijlagen.nsf/web/Bijlagen_Home_NI

                BIJLAGEN
                                                                          Identificatie van het formulier
 Laatste versie
 9 Functienummers
 10 Codificatie vergoedingen
                                                                                        12 Identificatie vh risico
 13 Code aard van de dag
  14 Afwijkende gebeurtenissen
                                                                                                       Begindatum geldigheid Einddatum geldigheid
  15 Betrokken voorwerpen
                                              Code
                                                          Omschrijving
                                                 AADD501 Aanvraag aangifte
                                                                                                      01-01-1900
                                                                                                                                  01-01-9999
                                                 ACRF001 Ontvangstbewijs
                                                                                                      01-01-1900
                                                                                                                                  01-01-9999
                                                 AOAT001 Arbeidsongevallen scenario 1 - aangifte van een
                                                                                                      01-01-1900
                                                                                                                                  01-01-9999
 Gewijzigde bijlagen
                                                 AOAT002 Arbeidsongevallen scenario 2 - maandelijkse
                                                                                                      01-01-1900
                                                                                                                                  01-01-9999
                                                          rapportering
                                                                                                                                  01-01-9999
 Per status
                                                 AOAT003 Arbeidsongevallen scenario 3 - mededeling van
                                                                                                      01-01-1900
                                                                                                                                  01-01-9999
                                                 AVDDTUP Wijzigende aangifte met betrekking tot een aangifte 01-01-1900
 Overzicht
                                                 AVWDDT Aangifte van werken
                                                                                                                                  01-01-9999
                                                 BEWARE Notificatie die de boekhoudkundige informatie bevat 01-01-1900
                                                                                                                                  01-01-9999
                                                          in verband met de wijzigingen van de multifunctionele
```

```
kg:IdentificationShape a sh:PropertyShape;
    rdfs:comment "Property Shape for Identification (code)";
    sh:path ont:Identification; sh:datatype xsd:string;
    sh:maxLength 7; sh:minCount 1; sh:maxCount 1;
    # Geldige waarden komen van bijlage 1
    sh:in ( "AADD501" "ACRF001" ... "ZIMA006" ) ;
    # De geldige waarden hebben ook een geldigheidsdatum,
    # we demonstreren dit voor 1 entiteit
    sh:sparql [
      sh:message "Date not within valid interval of document type";
      sh:prefixes <> ;
        sh:select """
          SELECT $this (ont:Identification AS ?path) ?value
          WHERE
            $this ont: Identification ?value .
            $this ont:FormCreationDate ?date .
            ?idtype a ont:IdentificationType ; rdfs:label ?value ;
                    ont:validFrom ?from ; ont:validTo ?to .
            FILTER ( !( ?date >= ?from && ?date <= ?to ) )
          3 " " " ;
```



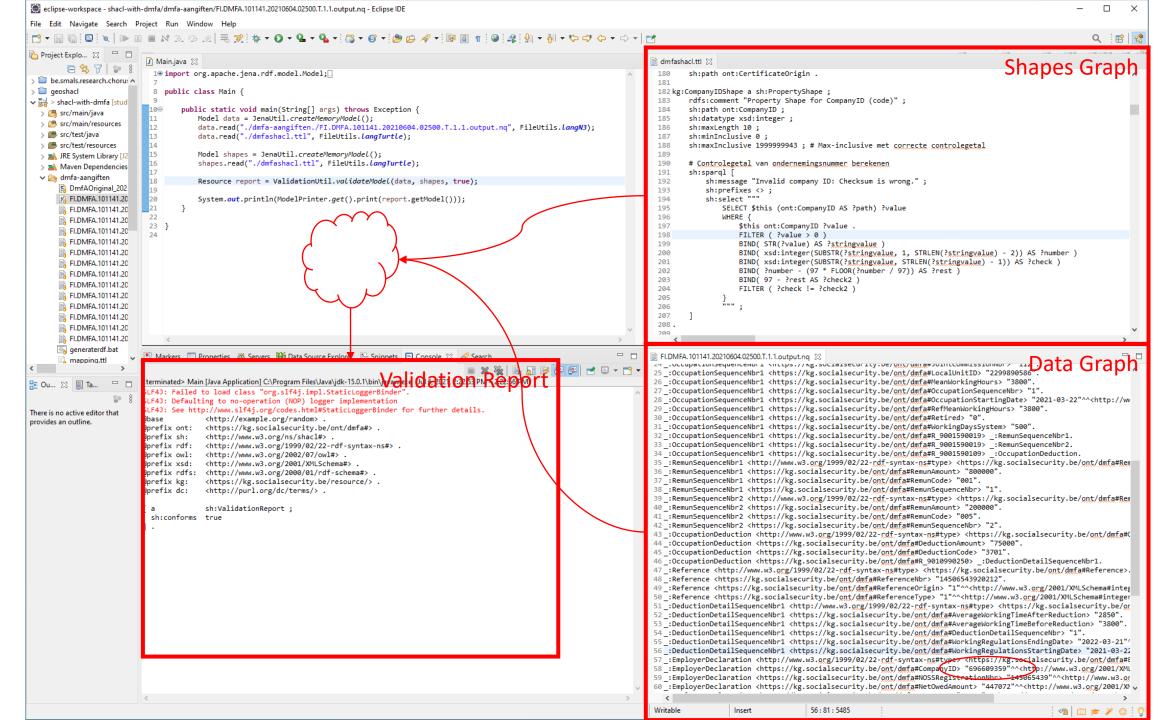
# A small demonstration

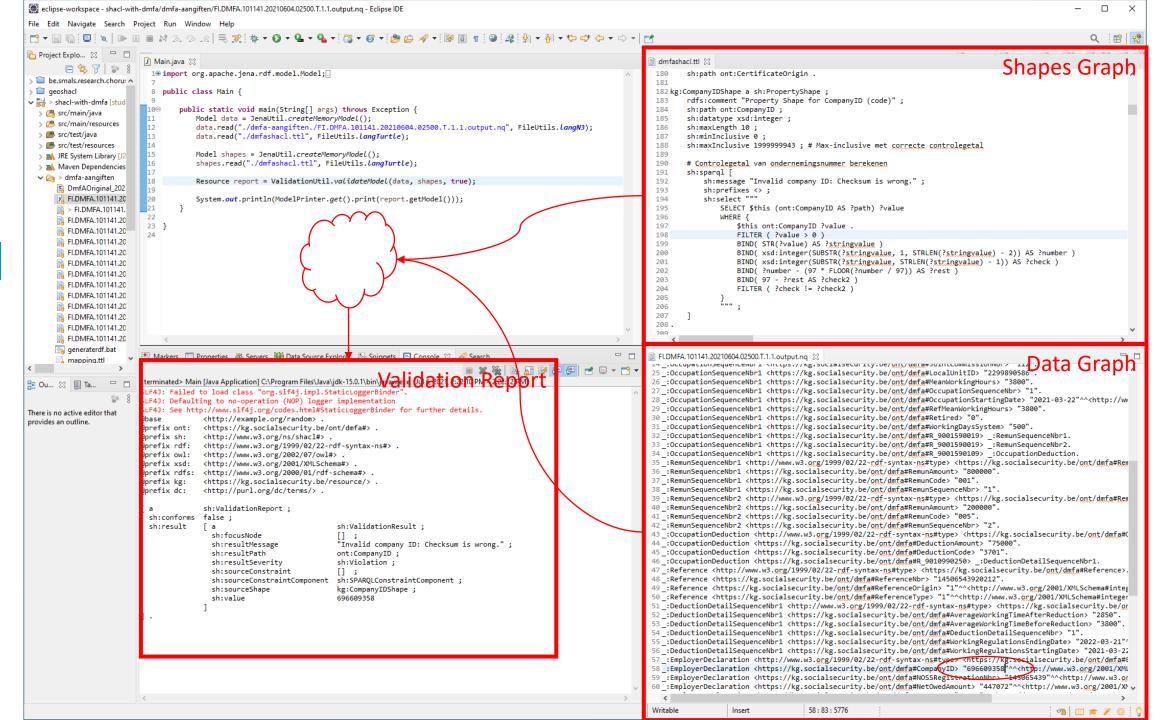
### Demonstratie

- By means of the TopBraid library
- Data graph = DmfA declaration in RDF + extra data
- Shapes graph

Focus on calculation control number







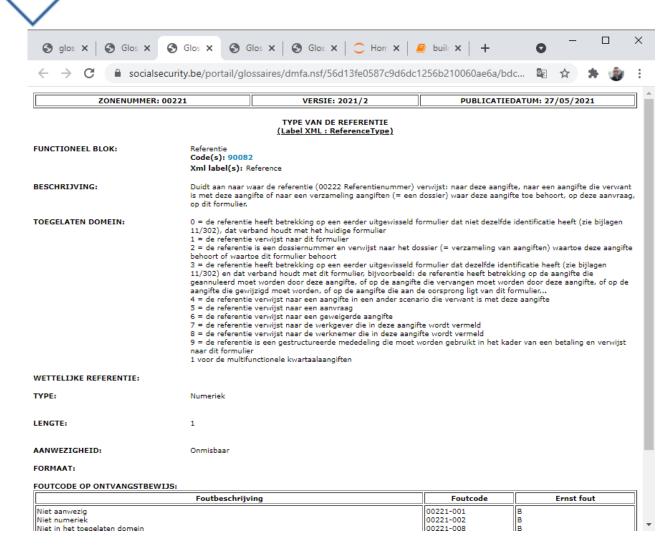
```
kg:CompanyIDShape a sh:PropertyShape;
   rdfs:comment "Property Shape for CompanyID (code)";
   sh:path ont:CompanyID ; sh:datatype xsd:integer ; sh:maxLength 10 ;
   sh:minInclusive 0; sh:maxInclusive 1999999943; # Max-inclusive met correcte controlegetal
   # Controlegetal van ondernemingsnummer berekenen (ofwel 0, ofwel correct volgens berekening
   sh:sparql [
     sh:message "Invalid company ID: Checksum is wrong.";
     sh:prefixes <> ;
     sh:select """
       SELECT $this (ont:CompanyID AS ?path) ?value
       WHERE {
         $this ont:CompanyID ?value .
         FILTER ( ?value > 0 )
         BIND( STR(?value) AS ?stringvalue )
         BIND( xsd:integer(SUBSTR(?stringvalue, 1, STRLEN(?stringvalue) - 2)) AS ?number )
         BIND( xsd:integer(SUBSTR(?stringvalue, STRLEN(?stringvalue) - 1)) AS ?check)
         BIND( ?number - (97 * FLOOR(?number / 97)) AS ?rest )
         BIND( 97 - ?rest AS ?check2 )
         FILTER ( ?check != ?check2 )
       } """;
```

Check company numbers! This rule can later be generalized so that it can be reused within an organization.

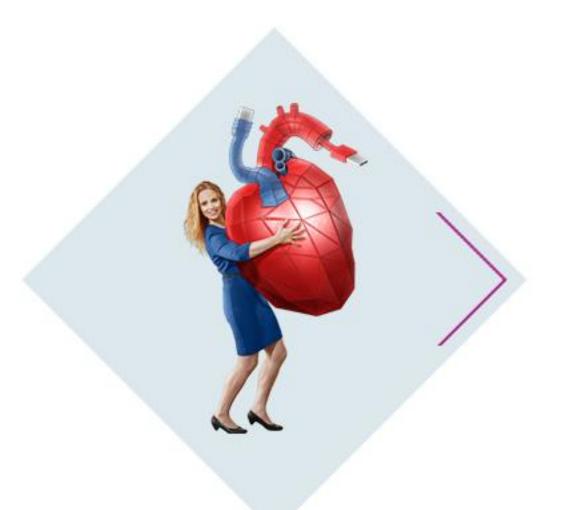


### Into the future?

- Some descriptions refer not only to external data (see example Appendix 11), but also to past declarations.
- This experiment involves local validation. Can we validate globally? In theory, yes, but then we have to store all declarations in a graph.







# Summary

# Summary

- Validating DmfA declarations as RDF is feasible
  - The generation of SHACL must be "aligned" with the generation of XSD
  - Sharing SHACL rules, like sharing the XSD, would allow third parties to validate their data locally.

SHACL is even more expressive than XSD

- SHACL and the vocabulary is compatible with JSON-LD (a JSON serialization of RDF), which is only in more contemporary Web environments (cfr. Flanders)
- SHACL can even validate declarations against a KG (future?)

