

PUBLIC MULTILINGUAL KNOWLEDGE MANAGEMENT INFRASTRUCTURE FOR THE DIGITAL SINGLE MARKET (2016.16)

IDENTIFICATION OF THE ACTION

| | |
|---------------------|--|
| Type of Activity | Common services, common frameworks |
| Service in charge | Publications Office of the European Union |
| Associated Services | DG Connect DG DIGIT DG DGT European Parliament, DG TRAD, Terminology Coordination Centre de Traduction |

EXECUTIVE SUMMARY

In their open letter to the European Commission the European Language (Technology) Community claims: Europe's Digital Single Market must be multilingual!¹ EU cross-border online services represent only 4% of the global Digital Market and only 7% of small and medium sized enterprises (SMEs) in the EU are actually selling cross-border.² Providing support for the EU economy and in particular to SMEs to overcome the language barriers will help to unlock the e-Commerce potential within the EU.

The objective of this action is to support enterprises and in particular the language technology industry with the implementation of the necessary multilingual tools and features in order to improve cross border accessibility of e-Commerce solutions. It will also be an input to the CEF Automated Translation Platform, a common building block implemented through the CEF programme to be used by European cross-border online public services.

The public multilingual knowledge infrastructure will be governed by a representative subgroup of stakeholders of the final system.

In this context multilingual tools and features refer to capabilities such as machine translation, localisation and multilingual search. The public multilingual knowledge infrastructure should reduce the investments of enterprises for the creation of their individual knowledge management systems by providing an agreed, open, reliable and persistent public core knowledge management system. This would also create space for innovation instead of wasting resources for redundant activities.

Public administrations and public entities within the EU will largely benefit from this initiative, in particular regarding the internationalisation of their e-Services. They will be able to share and to valorise existing taxonomies/terminologies and to extend their multilingual capabilities. This will also help to increase the interoperability between public administrations within EU in general.

¹ See <https://ec.europa.eu/futurium/en/content/european-language-and-language-technology-community-europes-digital-single-market-must-be>

² See http://europa.eu/rapid/attachment/IP-15-4653/en/Digital_Single_Market_Factsheet_20150325.pdf

To realise the public multilingual knowledge infrastructure the following aspects need to be addressed:

- Implementation of a technical infrastructure to expose existing multilingual taxonomies/terminologies in a standardised way based on semantic technology and Semantic Web standards;
- Implementation of existing alignments between terminologies and creation of further alignments and relations in order to enable interoperability;
- Creation and maintenance of meaningful supplements, i.e. of terms and relations that complete the coverage of the multilingual knowledge infrastructure and facilitate interoperability;
- Set-up of a community and a governance structure to extend systematically the coverage of the core infrastructure by the integration of supplementary public multilingual taxonomies/terminologies.

The cornerstone of the public multilingual knowledge infrastructure will be EuroVoc, the multilingual, multidisciplinary thesaurus covering the activities of the EU, which is managed by the Publications Office. The project will also benefit from already existing alignments of EuroVoc with other thesauri (Agrovoc, Eclas, Gemet, Anubis and Inspire³).

In addition, it should be investigated how the publication of the information as Linked Open Data (LOD) could be enhanced by introducing lexical semantic links between semantically equivalent and similar entities in an automatized way.

OBJECTIVES

The objectives of the proposed activity are:

- To provide an agreed, reliable, persistent and extensible public multilingual terminology platform for multiple purposes and for multiple stakeholders composed of open public multilingual resources.
- To increase the interoperability of existing multilingual terminologies, in particular by aligning and linking them with other existing terminologies. Linking will enable at the same time specialisation (for example by linking a concept of a more general taxonomy/terminology with the corresponding concept of a domain specific taxonomy/terminology) and broadening (for example by linking similar concepts at the same level of granularity).
- To establish the initial governance structure to support and to supervise the execution of the project as well as the implementation, the management and the evolution of the final system. Synergy with the governance structure, which is being set-up for the CEF.AT platform (see <http://www.lr-coordination.eu/anchor-points>), will be established.
- To contribute to the further standardisation of data models for thesauri and lexical knowledge database representations using latest semantic technologies.
- To support the contributing institutions with the transformation of their resources into adopted semantic format of the platform.
- To further develop the LOD capabilities of the system, in particular by enhancing the creation of semantic links between similar and related concepts.

³ See <https://open-data.europa.eu/en/data/dataset/eurovoc>

SCOPE

The objective of the proposed activity in the scope of the ISA² programme is to verify the feasibility of the approach and to prepare the technical and the organisation aspects for the definitive and permanent implementation of an open public multilingual knowledge infrastructure managed by the EU Institutions.

Nevertheless, it will deliver already results, which could be used and applied by public administrations and bodies of Member States and EU Institutions independent from the public multilingual knowledge infrastructure project itself.

In scope

1. Adoption of a standard representation for multilingual terminologies (candidates include SKOS⁴, Lemon⁵...)
2. Definition of a core data model based on the standard representation in order to facilitate the interoperability between different terminologies, i.e. through a shared set of metadata, and to harmonise the representation of the data, which will be made available through the platform
3. Specification of the technical architecture of the public multilingual knowledge infrastructure and the necessary services to access and to manage the system
4. Proof of concept – implementation of a first operational release of the system to demonstrate the core services of the system
5. Set-up/adoption of an adequate initial governance structure
6. Definition of an iterative implementation strategy, i.e. the specifications and the roadmap for the extension of the initial release of the system into a public service, which will be managed, further developed and maintain by the EU Institutions and governed by all contributors, i.e. in particular public entities within in EU. The result of the proof of concept should be reused for the implementation of the final system.
7. Feasibility study in order to analyse and to test the creation of lexical semantic links between semantically equivalent and similar entities in an automatized way.

Out of scope

1. Implementation of the definitive, permanent platform, made available as a public service and free of charge.
2. Management, further development and maintenance of the definitive system.

ACTION PRIORITY

The creation of a Public Multilingual Knowledge Infrastructure contributes to the overcoming of language barriers, in particular in the context of the implementation of a digital single market. It should also help to reduce the investments of the different stakeholders in cross-border e-commerce solutions and multi-lingual eGovernment solutions and to enhance the linguistic quality of the solutions.

⁴ See <http://www.w3.org/2004/02/skos/>

⁵ See <http://lemon-model.net/index.php>

Contribution to the interoperability landscape

The contribution of the action to the interoperability landscape, measured by the importance and necessity of the action to complete the interoperability landscape across the Union

| Question | Answer |
|---|--|
| <p><i>Does the proposal directly contribute to implementing the European Interoperability Strategy, the European Interoperability Framework, or other EU policies with interoperability requirements, or needed cross-border or cross-sector interoperability initiatives? If yes, please indicate the EU initiative / policy and the nature of contribution.</i></p> | <p>Yes. The proposal meets the recommendations included in the European Interoperability Strategy (EIS)⁶. In particular the adherence to specific standards for describing language resources and the creation of an interoperability platform to manage them comply with the main approaches and “clusters” of the EIS (reusability of the solutions, interoperability service architecture in the EU multilingual context, implication of ICT on new EU legislation, as well as promotion of the awareness on the maturity level and of the shareability of the public administration services). Similarly, the proposal meets the recommendations and principles of the European Interoperability Framework (EIF)⁷, in particular as regards multilingualism, accessibility, administrative simplification, transparency, reusability of the solutions. The creation of a public multilingual knowledge infrastructure will allow EU public administrations to create services that can be accessible and shareable independently from the language actually used, as well as the SMEs to sell goods and service cross-border in a digital single market.</p> |
| <p><i>Does the proposal fulfil an interoperability need for which no other alternative solution is available?</i></p> | <p>Yes. This action represents a tremendous opportunity to harmonize the different language resources managed by EU institutions (for example Eurovoc, IATE, Glossaries searchable on GlossaryLinks of the DG TRAD, etc.), as well as the national resources managed by Member States, and make them interoperable.</p> |

⁶ COM(2010) 744 final Annex 1, http://ec.europa.eu/isa/documents/isa_annex_i_eis_en.pdf

⁷ COM(2010) 744 final Annex 2, http://ec.europa.eu/isa/documents/isa_annex_ii_eif_en.pdf

1.1.1.1 Cross-sector

The scope of the action, measured by its horizontal impact, once completed, across the sectors concerned

| Question | Answer |
|---|--|
| <p>Will the proposal, once completed be useful, from the interoperability point of view, and utilised in two (2) or more EU policy areas? If yes, which are those?</p> | <p>This action aims at establishing multilingual interoperability of language resources; therefore it will promote multilingual interoperability services, as cross-collection and cross-language information retrieval, as well as translation services.</p> <p>It will contribute therefore to facilitate the creation of a Digital Single Market in the EU, which represents one of the main priorities of the European Commission. In particular it addresses all the three policy areas of such priority:</p> <ul style="list-style-type: none"> - Better online access to digital goods and services - An environment where digital networks and services can prosper - Digital as a driver for growth. |
| <p>For proposals or their parts already in operational phase: have they been utilised in two (2) or more EU policy areas? Which are they?</p> | <p>This proposal is not in operational phase yet.</p> |

Cross-border

The geographical reach of the action, measured by the number of Member States and of European public administrations involved.

| Question | Answer |
|--|--|
| <p>Will the proposal, once completed be useful, from the interoperability point of view, and used by public administrations of three (3) or more EU Members States?</p> | <p>By guaranteeing interoperability of language resources in all the 24 official languages of the EU, this proposal has the potential of improving the service interoperability of public administrations of all EU Member States, candidate countries or EFTA States.</p> |
| <p>For proposals or their parts already in operational phase: have they been utilised by public administrations of three (3) or more EU Members States?</p> | <p>This proposal is not in operational phase yet.</p> |

Urgency

The urgency of the action, measured by its potential impact, taking into account the lack of other funding sources

| Question | Answer |
|--|--|
| <i>Is your action urgent? Is its implementation foreseen in an EU policy as priority, or in EU legislation?</i> | <p>The outcomes of this action can greatly improve the accessibility of EU and Member States' legislation and related information systems by promoting the interoperability of the language resources used for automatic translation, as well as multilingual classification and indexing. Moreover, it will promote e-commerce solutions and related services which will rely on an agreed, authentic and persistent set of multilingual terminologies. This action is in particular foreseen in the framework of the creation of a European multilingual Digital Single Market, which is one of the priorities of the European Commission, aimed at supporting the EU economy (in particular the SMEs) to overcome the language barriers in order to unlock the e-Commerce potential within the EU.</p> <p>A prompt implementation of such proposal will have direct beneficial impacts on the addressed fields.</p> |
| <i>Does the ISA² scope and financial capacity better fit for the implementation of the proposal as opposed to other identified and currently available sources?</i> | <p>Yes. The proposal is specifically addressed to improve the <i>interoperability of language resources</i> and the <i>services for public administration and SMEs</i>. For both these reasons, ISA² fits to it better than other actions.</p> |

Reusability of action outputs

The re-usability of the action, measured by the extent to which its results can be re-used

Can the results of the proposal be re-used by a critical part of their target user base, as identified by the proposal maker? For proposals or their parts already in operational phase: have they been re-used by a critical part of their target user base?

| | |
|--|--|
| Name of reusable solution | Core data model for Multilingual taxonomies/terminologies |
| Description | <p>Formal definition of the core data model for multilingual taxonomies/terminologies and its necessary extensions that will be implemented by the public multilingual knowledge infrastructure.</p> <p>The approach should be flexible in the way that data providers would be able to define private extensions, which would allow the upload of supplementary data that is available on their side and that could be useful for re-users.</p> <p>The aspects "provenance" and "license" have also to be taken into account.</p> |
| Reference | PUB_MUL_TERM_FORMAT |
| Target release date / Status | Q2/2017 |
| Critical part of target user base | n/a |
| For solutions already in operational phase - actual reuse level (as compared to the defined critical part) | Not in operational phase |

| | |
|--|---|
| Name of reusable solution | Semantic links |
| Description | <p>Feasibility study and prototype in order to explore the possibilities to enhance the semantic capabilities of the platform, in particular regarding the creation of lexical semantic links between semantically equivalent and similar entities in an automatized way.</p> |
| Reference | PUB_MUL_TERM_SEMANTIC |
| Target release date / Status | Q4/2017 |
| Critical part of target user base | n/a |
| For solutions already in operational phase - actual reuse level (as compared to the defined critical part) | Not in operational phase |

| | |
|---------------------------|--|
| Name of reusable solution | First release of the system |
| Description | <p>Implementation of a first release of the system (repository and core services), which should be considered first of all as an operational proof of concept to demonstrate the core services of the platform and which will be reused to build the final system.</p> |
| Reference | PUB_MUL_TERM_POC |

| | |
|--|--------------------------|
| Target release date / Status | Q3/2018 |
| Critical part of target user base | n/a |
| For solutions already in operational phase - actual reuse level (as compared to the defined critical part) | Not in operational phase |

Level of reuse by the proposal

The re-use by the action of existing common frameworks and elements of interoperability solutions.

| Question | Answer |
|---|--|
| Does the proposal intend to make use of any ISA ² , ISA or other relevant interoperability solution(s)? Which ones? | The proposal will make use of VocBench 3, developed in the context of the ISA ² 2016. VocBench 3 represents a direct evolution of VocBench 2, originally developed by the Food and Agriculture Organization (FAO) of the United Nations for specifically managing their thesaurus Agrovoc and later evolved into a general purpose SKOS editor, adopted, among others, by the Publications Office for the maintenance of EuroVoc. |
| For proposals or their parts already in operational phase: has the action reused existing interoperability solutions? If yes, which ones? | Not in operational phase |

Interlinked

The link of the action with Union initiatives to be measured by the collaboration and contribution level of the action to Union initiatives such as the DSM.

| Question | Answer |
|--|---|
| Does the proposal directly contribute to at least one of the Union's high political priorities such as the DSM? If yes, which ones? What is the level of contribution? | This proposal contributes in particular to Digital Single Market (DSM) priority (cfr. 1.1.5.2). |

PROBLEM STATEMENT

"Linguistic diversity is and must remain a cornerstone and treasured cultural asset of Europe. However, the language barriers created by our 24 official EU languages cause the European market to be fragmented and to fall short of its economic potential. Almost half of European citizens never shop online in languages other than their native tongue, access to public e-services is usually restricted to national languages and the richness of EU educational and cultural content is confined within linguistic communities. European SME's are at particular disadvantage, because the cost of providing services in multiple languages is prohibitive and has a negative impact on their competitiveness."⁸

This challenge needs to be addressed and a public multilingual knowledge infrastructure will contribute to reduce and to secure the investments of the different stakeholders in cross-border e-commerce solutions and related services because part of their implementation could rely on an agreed, authentic and persistent set of multilingual terminology.

Because the contributions for public multilingual knowledge infrastructure will come from different stakeholders (essentially public administrations and bodies of EU Member States, EU Institutions and international organisations) the challenge is to build a system, which empowers the stakeholders to manage the development and evolution of their taxonomies/terminologies on an individual base, but at the same time enables interoperability through alignment and linking.

The only constraint should be that stakeholders have either to adopt the core data model proposed by the public multilingual knowledge infrastructure or, at least, have to be able to perform the necessary transformations to provide new releases in compliance with the core data model. Ideally, new releases have to be provided in a way that they can be integrated largely automatically.

EXPECTED BENEFICIARIES AND ANTICIPATED BENEFITS

| Beneficiaries | Anticipated benefits |
|---------------------------------|--|
| EU economy | Many studies have already been conducted to evaluate the possible economic impact of an increase in cross-border e-commerce between EU Member States. The creation of a real EU digital single market has become a priority of the Commission. The initiative will provide a contribution on the technological level. It will help to reduce the localisation effort for e-commerce platforms, enhance the quality of the domain specific terminology and improve their harmonisation. It will also facilitate the implementation of multilingual search capabilities. |
| EU language technology industry | Cost reduction and faster implementation of services related to cross-border e-commerce (machine-translation, localisation software, cross-language search solutions...). This will also increase the usability and searchability of resources for the creation of new, innovate services. |
| EU Member States | Will benefit in the context of the internationalisation of their e-government services for example to facilitate foreign investments in the local market. |

⁸ See <https://ec.europa.eu/futurium/en/content/european-language-and-language-technology-community-europes-digital-single-market-must-be>

| | |
|-----------------|---|
| | Will be able to improve interoperability with other Member States and/or public entities based on shared or aligned taxonomies/terminologies. |
| EU Institutions | Valorisation of existing multilingual taxonomies/terminologies, spin-offs for EU translation services and other multilingual services. It will help to increase the interoperability of multilingual LOD, which are made available by the EU Institutions. |

EXPECTED MAJOR OUTPUTS

| | |
|------------------------------|--|
| Output name | Technical architecture |
| Description | Technical design of the public multilingual knowledge infrastructure architecture including definition of all related services (ingestion of and access to data (including search), management of the infrastructure itself...). |
| Reference | PUB_MUL_TERM_ARCHITECTURE |
| Target release date / Status | Q3/2017 |

| | |
|------------------------------|--|
| Output name | Governance structure |
| Description | Proposal for an adequate governance structure for the supervision of the public multilingual knowledge infrastructure. |
| Reference | PUB_MUL_TERM_ARCHITECTURE |
| Target release date / Status | Q2/2017 |

| | |
|------------------------------|---|
| Output name | Implementation strategy |
| Description | Proposal of an iterative implementation strategy in order to prepare the political decision about whether the EU institutions will support the implementation of the public multilingual knowledge infrastructure and, if yes, how the system should be managed and financed. |
| Reference | PUB_MUL_TERM_STRATEGY |
| Target release date / Status | Q1/2019 |

| | |
|------------------------------|--|
| Output name | Community building |
| Description | Proposal for implementation and organisation of a community (contributors, users...) to drive the further evolution of the system and of the language resources. |
| Reference | PUB_MUL_TERM_COMMUNITY |
| Target release date / Status | Q2/2019 |

ORGANISATIONAL APPROACH

Expected stakeholders and their representatives

| Stakeholders | Representatives |
|------------------------------|---|
| EU Institutions | European Parliament DG TRAD, Terminology Coordination Commission DIGIT DG DGT DG CONNECT Publications Office of the EU Translation Centre for the Bodies of the European Union |
| International organisations | FAO |
| Member states | |
| Language technology industry | For example the companies represented by the LT innovate Association |
| Research community | Universities and research institutions that are active in this area |

Identified user groups

- Member States (public administrations involved in "internationalisation" and "eGovernment" initiatives)
- Implementers of eGovernment solutions
- European Institutions and bodies
- Language Technology Industry and their customers
- Citizens
- Candidate countries, EFTA and other countries(public administrations involved in "internationalisation" and "eGovernment" initiatives)

Communication plan

The following table presents a first rough idea of a communication plan based of the different beneficiaries/interest groups, which have been identified in a first phase.

The existing platforms of the ISA programme in the domain of language technology will be taken into account for the set-up of an adequate communication platform.

| Beneficiaries | Communication channel | Frequency |
|---------------------------------|--|--|
| EU economy | Web (information about the activity on the ISA ² website, publicity on the Publications Office and other EU Institutions websites) | Regular updates during the lifetime of the project. |
| EU language technology industry | Web (information about the activity on the ISA ² website, publicity on the Publications Office and other EU Institutions websites) Conferences (delivery of presentations) | Regular updates during the lifetime of the project. 1 to 5 conferences per year |
| Member States | Web (information about the activity on the ISA ² website, publicity on the Publications Office and other EU Institutions websites) Workshops (organisation of dedicated workshops with interested member states) | Regular updates during the lifetime of the project. 1 to 3 workshops per year |
| EU Institutions | Meetings Workshop (organisation of dedicated workshops with interested services) | Regular meetings of the EU institutional stakeholders 1 to 3 workshops per year |
| Terminology community | Conferences (delivery of presentations) | 1 to 3 conferences per year |
| Semantic Web community | Conferences (delivery of presentations: SEMIC, dedicated conferences...) | 1 to 3 conferences per year |

Governance approach

The implementation of a governance body is needed at different levels.

In the first phase the work to be done in the scope of the ISA² programme needs to be governed by a representative subset of the stakeholders of the final system. This group is considered as the implementation of the initial government structure. It should rely on governance structure, which is being set-up for the CEF.AT platform (see <http://www.lr-coordination.eu/anchor-points>).

If the implementation of the system has been decided, the governance structure has to be adapted to be able to support and to supervise the management and the further evolution of real production system.

TECHNICAL APPROACH AND CURRENT STATUS

Aspects to be considered:

- Management of multilingual taxonomies/terminologies

The data structure for the public multilingual knowledge infrastructure will be defined by a core data model, which will be composed of two parts: a mandatory part (core metadata), which has to be respected by all data providers and optional part (private extensions) to allow the publication of additional data, which exist for a particular dataset and which are not covered by the core data model, if it represents an added value for the users of the system. I.e. additional data could be stored by the system but will eventually not be fully supported by the common services offered by the system.

All individual concepts have to be represented in the adopted (semantic) format. Each individual concept will be identified by a unique persistent URI.

The reuse/adoption of existing software components will be encouraged.

- Distributed infrastructure

The public multilingual knowledge infrastructure should be implemented as a distributed network of RDF triple stores in order to guarantee a maximum of availability of the system.

Adequate management capabilities are needed to guarantee the consistency of the data.

APIs and online access should be implemented in a way that the technical implementation is hidden. The user works on a virtual system, which is composed by a set of federated RDF triple stores, physically hosted in different locations.

- Management of the system

Also the management services of the system should be implemented in a way that the technical implementation is hidden.

There will be different groups of services:

- Maintenance of data

 - Ingestion of new data sets (including validation processes)

 - Update of existing data sets

 - Management of supplementary concepts, i.e. concepts that only exist on the level of the public multilingual infrastructure (create, update, delete)

 - Search and visualisation

- Maintenance of data structure (core data model and extensions, relations, alignments...)

- Management of the platform itself

 - Administration interface (monitoring of services, configuration, user management (for contributors and administrators))

COSTS AND MILESTONES

Breakdown of anticipated costs and related milestones

| Phase: Initiation Planning Execution Closing/Final evaluation | Description of milestones reached or to be reached | Anticipated Allocations (KEUR) | Budget line ISA/ others (specify) | Start date (QX/YYYY) | End date (QX/YYYY) |
|--|---|--------------------------------------|---|-------------------------|-----------------------|
| Inception | Project organisation has been set-up | 60 | | Q3/2016 | Q4/2016 |
| Execution | Standard representation has been adopted | 50 | | Q4/2016 | Q1/2017 |
| Execution | Core data model and a first set of extensions have been defined (including documentation) | 100 | | Q4/2016 | Q2/2017 |
| Execution | Technical architecture has been defined | 100 | | Q2/2017 | Q3/2017 |
| Execution | Proposal for an adequate government structure has been defined | 50 | | Q1/2017 | Q2/2017 |
| Execution | First release of the system (operational proof of concept) | 300 | | Q1/2018 | Q3/2018 |
| Execution | Proposal for the implementation strategy | 60 | | Q4/2018 | Q1/2019 |
| Execution | Creation of the community | 60 | | Q4/2018 | Q2/2019 |
| Execution | Feasibility study for the enhancement of the semantic capabilities of the platform | 144 | | Q2/2017 | Q4/2017 |
| | Total | 924 | | | |

Breakdown of ISA funding per budget year

| Budget Year | Phase | Anticipated allocations (in KEUR) | Executed budget (in KEUR) |
|-------------|---------------------|-----------------------------------|---------------------------|
| 2016 | Inception/execution | 160 | |
| 2017 | Execution | 344 | |
| 2018 | Execution | 360 | |
| 2019 | | 60 | |
| 2020 | | | |