# Spanish Tax Agency

# ITS 2.0 implementation experience in HTML5:

# www.agenciatributaria.es

**MultilingualWeb Workshop**
**Making the Multilingual Web Work**
**Rome, 12–13 March 2013**

**Spanish Tax Agency, IT department**

Román Díez González

Spanish Tax Agency

Pedro L. Díez-Orzas

Linguaserve

Linguaserve collaborators:

Giuseppe Deriard-Nolasco, Pablo Nieto Caride, Consuelo Aldana, Félix Fernández

# What are we talking about?

1. Introducing the Spanish Tax Agency

2. www.agenciatributaria.es in the MLW-LT project

3. Shifting to HTML5

4. Experience in ITS2.0 annotation:
   a. Automatic annotation of new ITS2.0 metadata
   b. Reusing custom tags for ITS2.0 metadata annotation
   c. Manual ITS2.0 annotation

5. Next steps and some proposals

Agencia Tributaria

# (1) Spanish Tax Agency

**Spain: General Indicators 2011**

″Spain is a country regionally structured into 17 autonomous communities and 2 autonomous cities with **5 co-official languages**

• Population : 47.190.493 inhabitants ( **12,2 % foreign residents**)

•**Spanish Tax Agency mission**

•Effective application of Spain's tax and customs structure

• Management of tax resources on behalf of other public administrations when ordered by Law or Agreement

•**Overall census of obliged taxpayers**

•Individual taxpayers:          46.509.231

• Companies:                      2.674.547

• Other organisations:          2.293.939

**Total taxpayers:**                **51.477.717**

GOBIERNO DE ESPAÑA          Agencia Tributaria

# What are we talking about?

1. Introducing the Spanish Tax Agency

2. www.agenciatributaria.es in the MLW-LT project

3. Shifting to HTML5

4. Experience in ITS2.0 annotation:
   a. Automatic annotation of new ITS2.0 metadata
   b. Reusing custom tags for ITS2.0 metadata annotation
   c. Manual ITS2.0 annotation

5. Next steps and some proposals based on experience
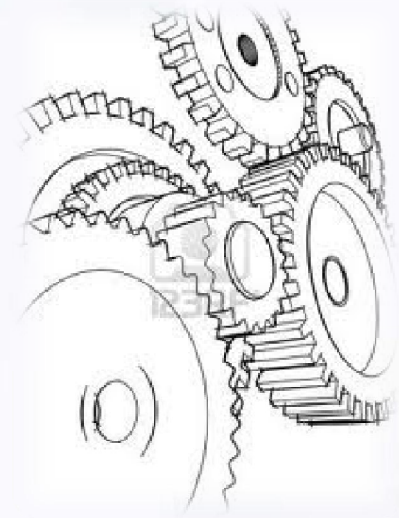
Agencia Tributaria

# (2) The Spanish Tax Agency in MLW-LT

- **www.agenciatributaria.es**, user in the "Online MT System" use case in the MultilingualWeb-LT (MLW-LT).

- The MLW-LT Working Group is administered by W3C and receives EC funding (LT-Web) through FP7 in the area of Language Technologies

# (2) The Spanish Tax Agency in MLW-LT

- Online MT System use case components:
  - Multilingual www.agenciatributaria.es
    (CMS: OpenText WEM)
  - HTML5
  - ITS 2.0
  - Real-time Multilingual Publication System
    - ATLAS (Linguaserve's Real Time Translation System)
    - Lucy Software MT (Rule-based Machine Translation)
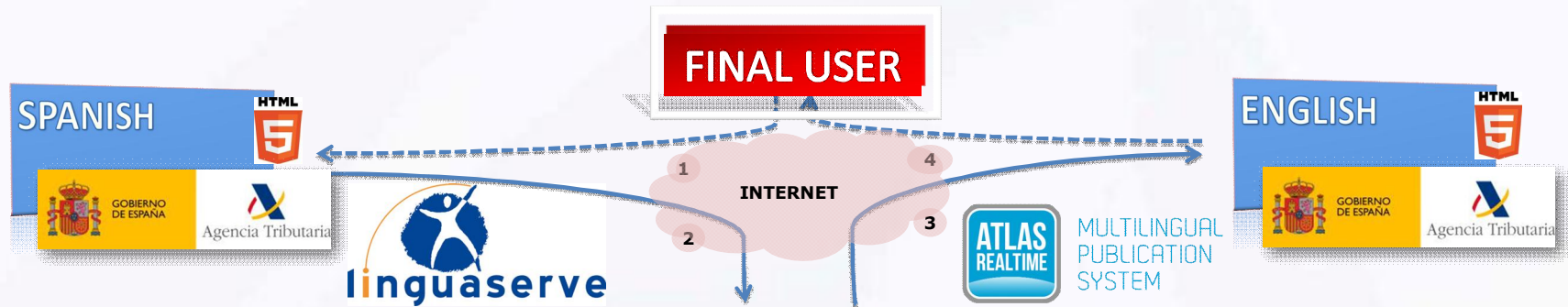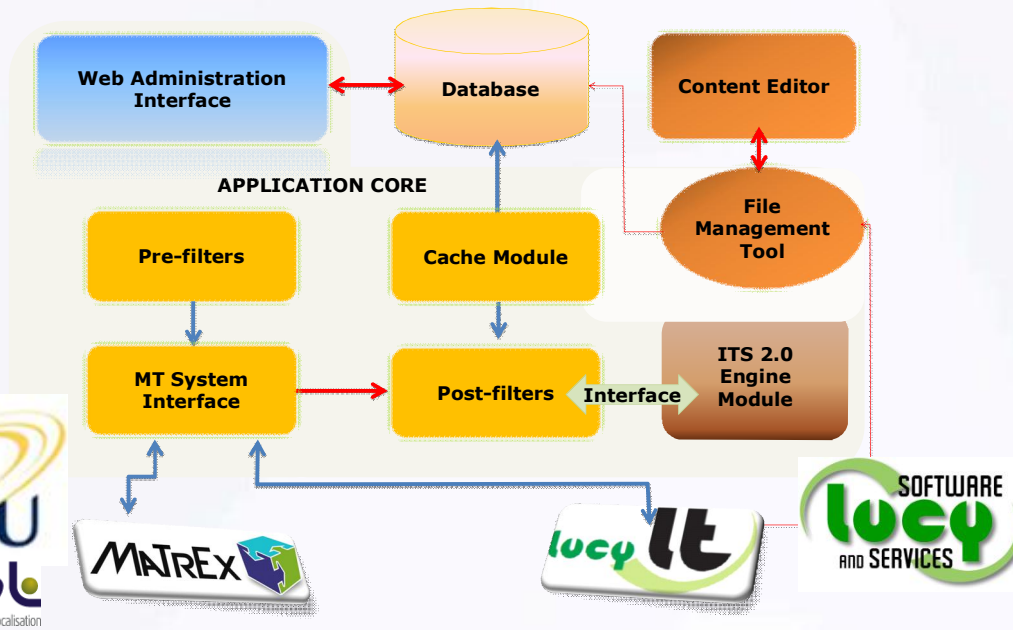    - MaTrEx from Dublin City University (Statistical Machine Translation)

Agencia Tributaria

# (2) Online MT System Use Case State

- RTMPS Implementation
  - Prototype 100% (ITS 2.0 definition from Dec 2012)
  - Showcase: preproduction demo (http://its2-aeat.linguaserve.net)
  - ITS 2.0 data categories: 6 (Translate, Localization Note, Language Information, Domain, Provenance, Localization Quality Issue)
- ES-EN total scope: 250 web pages. State:
  - Source language: 30% of target
  - Target language and Post-editing: 30% of target
- ES-FR, ES-DE total scope: 30 web pages. State:
  - Source language: 50% of target
  - Target language and Post-editing: 50% of target
- Testing: pending

MultilingualWeb-LT

# (2) Online MT System I18N



Please, see POSTER 4

# (2) MLW-LT Online MT SWOT

## Strengths

RTMPS highly reduces:

˝Translation costs (Quality on-demand)
   . MT + depending on % of post-editing cost reduction increases.

˝Management costs
˝Delivery time
˝Non-invasive technology

## Weaknesses

Viability dependent on :

˝Language combination
˝MT system output

˝Pre-editing and post-editing methodologies and tools (ITS 2.0 and HTML5 compliance)

## Opportunities

Profitability:

˝Websites with more than half a million words

˝Websites with a very high update frequency

## Threats

Control, performance and security:

˝The client might lose control of the translation → user¢s control with ITS 2.0

ÉReal-time performance
ÉSecurity level

# (2) ITS 2.0 benefits for the Spanish Tax Agency

- ITS 2.0 Increases user's control and automatic decision processes:
  - Translatability and language pair selection (Translate, Language information)
  - Specific terminology to apply (Domain)
  - Activation rules for post-editing (Localization Note)
  - Quality aspects reported to translation consumer or post-editor (Localization Quality Issue)
  - Post-editors judge quality of translation (MT Confidence)*
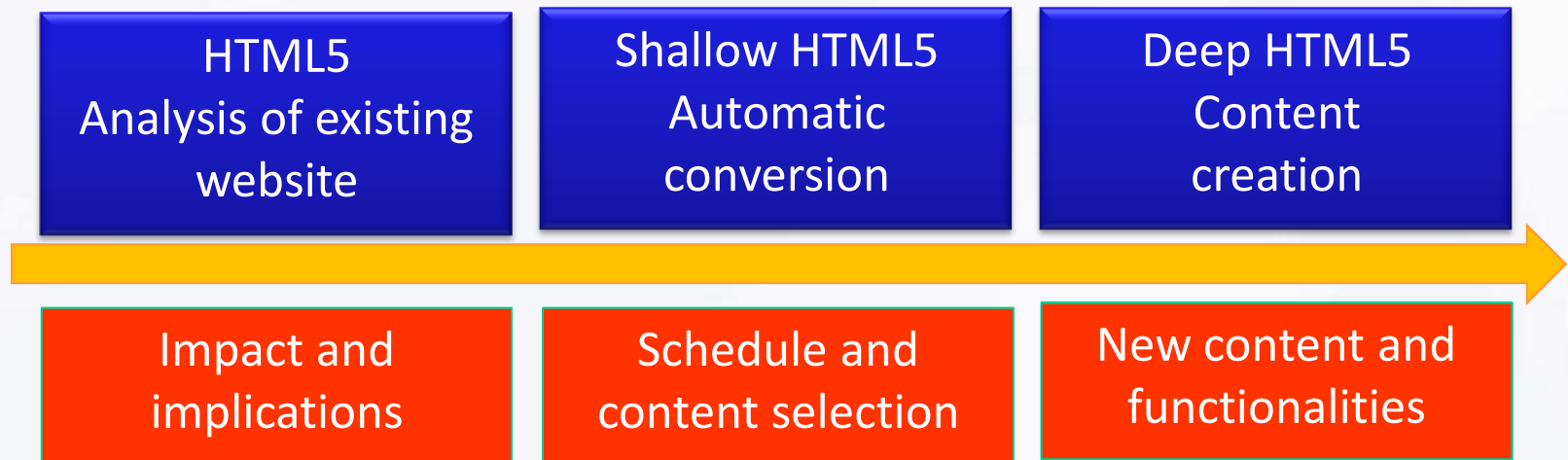  - Identification of agents (provenance)

# What are we talking about?

1. Introducing the Spanish Tax Agency

2. www.agenciatributaria.es in the MLW-LT project

3. Shifting to HTML5

4. Experience in ITS2.0 annotation:
   a. Automatic annotation of new ITS2.0 metadata
   b. Reusing custom tags for ITS2.0 metadata annotation
   c. Manual ITS2.0 annotation

5. Next steps and some proposals based on experience

# (3) Shifting to HTML5: Strategy

- Using ITS 2.0 requires HTML version 5 according to the current W3C specification.

| HTML5 Analysis of existing website | Shallow HTML5 Automatic conversion | Deep HTML5 Content creation |
|---|---|---|
| Impact and implications | Schedule and content selection | New content and functionalities |

# (3) Shifting to shallow HTML5: Modifications

- HTML5 DOCTYPE
- The language page (ISO 639-ISO 3166)
- Self-closed tags not allowed
- Head tags
- Erroneous nesting tags
- Attributes separated by spaces
- Non inclusion of presentation attributes in tags
- Header and body structure needed by tables

- HTML entities instead of special characters
- URLs cannot contain special characters
- ID attribute cannot contain spaces
- Required attributes (e.g. tag "object" must always have the attributes "data" and "type")
- Assessed attributes (e.g. "rel" attribute of tags "a" and "link" must be one from a closed list)

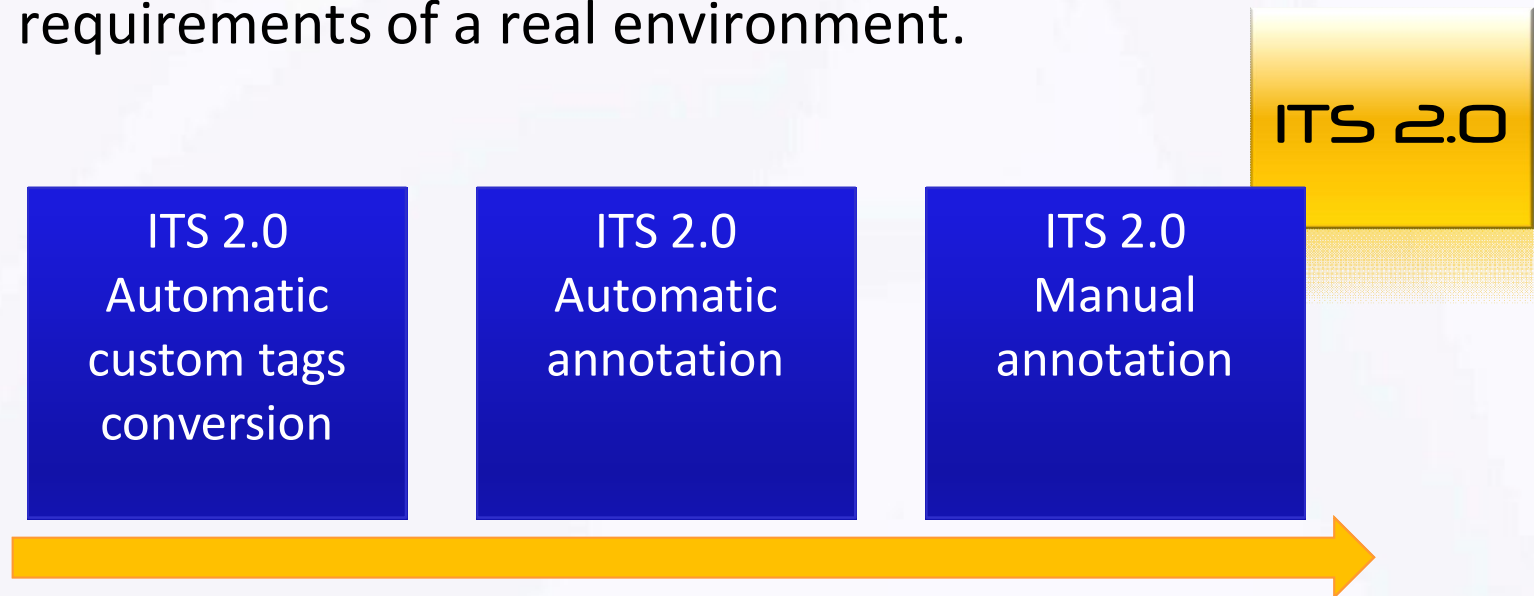# (3) Shifting to shallow HTML5: Obsolete attributes

| Tags | Impact |
|------|--------|
| **input** | Removed the alt attribute from any input tag that does not contain the attribute "type = 'image'" |
| **div** | Cannot define a "name" attribute in a "DIV" tag |
| **a** | Not allowed to define the attributes "name" and "title" in tag "a" |
| **embed and object** | Cannot define the attributes:<br>•"Applet" in the "embed" and "object" tags<br>•"Name" in the "embed" tag<br>•"Code", "archive", "classid", "codebase", "codetype", "state" and "standby" in the "object" tag |
| **table** | Not allowed to define the attributes "summary" and "border" in the "table" tag |
| **img** | Not allowed to define the attributes "name" and "border" in the "img" tag |
| **option** | Cannot define the attribute "name" in the "option" tag. |
| **param** | Not allowed to define the attributes "type" and "valuetype" in the "param" tag |
| **script** | Not allowed to define the attribute "lang" except in "JavaScript", it being case-insensitive in the tag "script" |
| **br** | Cannot define the attribute "clear" in the "br" tag |
| **background attribute** | No attribute is used to define the "background" in the tags "body", "table", "thead", "tbody", "tfoot", "tr", "td" and "th". |

# What are we talking about?

1. Introducing the Spanish Tax Agency

2. www.agenciatributaria.es in the MLW-LT project

3. Shifting to HTML5

4. Experience in ITS2.0 annotation:
   a. Automatic annotation of new ITS2.0 metadata
   b. Reusing custom tags for ITS2.0 metadata annotation
   c. Manual ITS2.0 annotation

5. Next steps and some proposals based on experience

# (4) ITS2.0 annotation experience

- Strategy adopted in order to annotate the content with ITS2.0 in an efficient and pragmatic way, considering the pressure and requirements of a real environment.

ITS 2.0

| ITS 2.0 Automatic custom tags conversion | ITS 2.0 Automatic annotation | ITS 2.0 Manual annotation |

# (4) Automatic ITS2.0 reuse of custom tags

É  Custom õno translateö tag already exists in the content and is automatically annotated as ITS 2.0 *Translate* data category:

*<li><!--ATLASP1NOTRAD--><a target="_blank" href="http://www.boe.es/diario_boe/txt.php?id=BOE-A-2011-20472">Orden EHA/3552/2011, de 19 de diciembre [í  ] <!--/ATLASP1NOTRAD--></li>*

*<li><a translate=önoö target="_blank" href="http://www.boe.es/diario_boe/txt.php?id=BOE-A-2011-20472">Orden EHA/3552/2011, de 19 de diciembre [í  ] </li>*
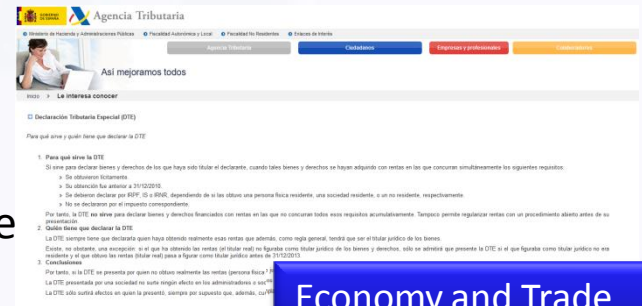
*Respecting the behaviour of the previous tag and the precedence rules of ITS:

   ó   Addition of ITS default rules for known translatable attributes:
   ó   <its:translateRule selector="//h:*/@title" translate="yes"/>
   ó   <its:translateRule selector="//h:*/@alt" translate="yes"/>

ITS 2.0

# (4) Automatic ITS 2.0 annotation: Domain

1. Extracting relevant domains based on the content.

2. Alignment of the domains with each web page.

3. Use of scripts and regular expressions to annotate the content.

**Economy and Trade**

4. Document processing:

   i. The selector points to the html root element, indicating that the domain applies to the whole **HTML** document (inheritance).

   ii. The **domainPointer** attribute indicates where the domain that applies to the selected content is ("**Economy and Trade**").

   iii. The **domainMapping** maps the domain "Economy and Trade" to "**ECON**", which will be sent as an understandable parameter to the MT System.

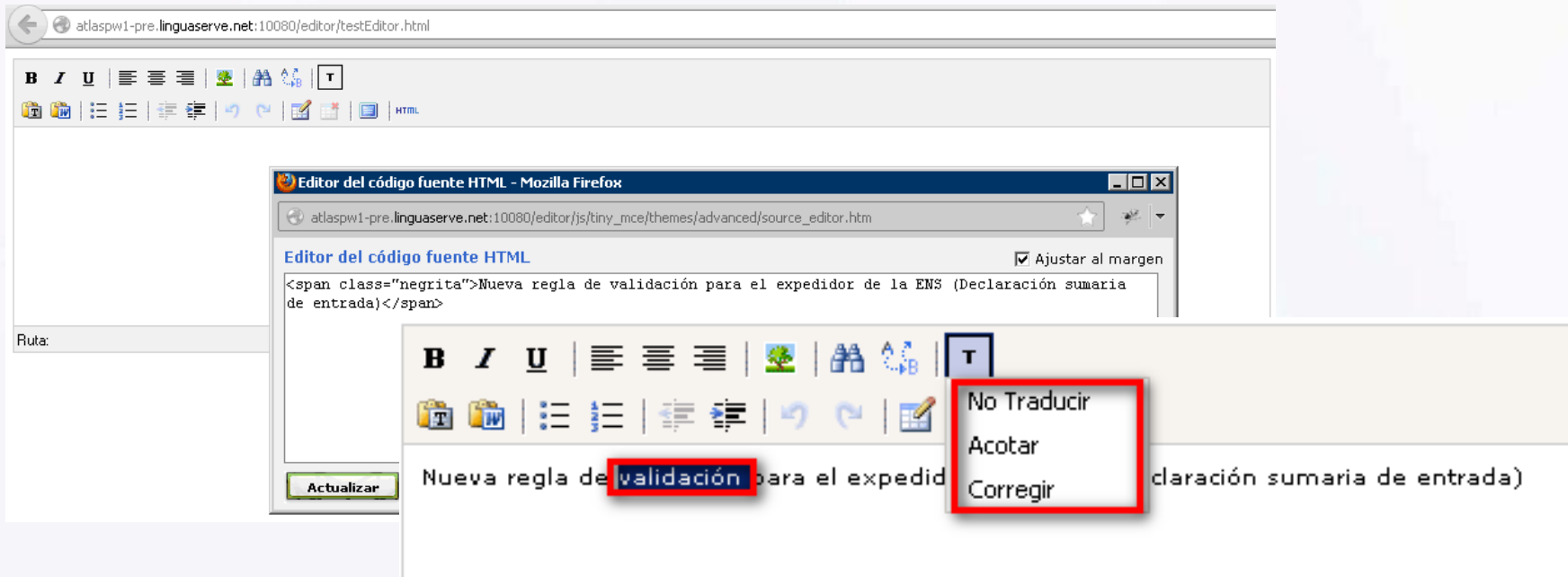```
<!DOCTYPE html>
<html lang="es">
<head>
<meta charset="utf-8">
<meta name="keywords" content="Economy and Trade"/>
[DOMAIN RULES]
</head>
<body>
[í  ]
</body>
</html>
```

```
<its:rules xmlns:its="http://www.w3.org/2005/11/its"
xmlns:h="http://www.w3.org/1999/xhtml" version="2.0">
<its:domainRule
selector="//h:html"
domainPointer="/html/head/meta[@name='keywords']/@contentõ
domainMapping="'Economy and Trade' ECON, 'Law and Legal
Science' LAW, -General Vocabulary' GV"/>
</its:rules>
```

MT System

# (4) Manual ITS2.0 annotation: Tool

- Quick and pragmatic approach:
  - New HTML Editor plugin created for the ITS 2.0 manual annotation for open source HTML Editor
  - User-friendly interface for the manual insertion of tags.

# (4) ITS 2.0 Manual annotation: Translate

- The author must only select the non-translatable element, click on the insertion icon (T) and click on the annotation type: **No Traducir**.

# (4) ITS 2.0 Manual annotation: Localization Notes

- Use of the annotation type **Acotar**:  The author inserts the annotation text into the box and the software will automatically create the tag.

É  The pull-down menu is used to choose the type of localization note. It can either be <u>description</u> (*descriptiva*) or <u>alert</u> (*alerta*).

<p>La disposición trigésima quinta de la Ley del <span its-loc-note="Stands for 'Impuesto sobre la Renta de las Personas Físicas ', use acronym in target language" its-loc-note-type="description">IRPF</span></p>

# (4) ITS 2.0 Manual annotation: Localization Quality Issue

- Use of the annotation type *Corregir*:  The author chooses a type of issue from a pull-down menu, inserts a comment into the box (*Comentario*), chooses a severity level between 0 and 100 (*Severidad*) and an optional link to a reference document (*documento de referencia*), and the software will automatically create the tag.

atlaspw1-pre.linguaserve.net:10080/editor/js/tiny_mce/plugins/translate/quality.htm

Tipo corrección  Terminología errónea

Comentario

Severidad (0-100)

documento de referencia

Enviar consulta

ITS 2.0

Online filing can be done by the interested party or by someone representing them. In both cases, an electronic certificate X.509.V3 issued by the <span its-loc-quality-issue-comment="**Has previously been translated as 'Royal Mint'. Please be consistent.**" its-loc-quality-issue-type="inconsistency" its-loc-quality-issue-severity="70">**National Coin and Stamp Factory**</span>.

# What are we talking about?

1. Introducing the Spanish Tax Agency

2. www.agenciatributaria.es in the MLW-LT project

3. Shifting to HTML5

4. Experience in ITS2.0 annotation:
   a. Automatic annotation of new ITS2.0 metadata
   b. Reusing custom tags for ITS2.0 metadata annotation
   c. Manual ITS2.0 annotation

5. Next steps and some proposals

# (5) Next steps and some proposals

- End of Online Translation System MLW-LT use case – June 2013

- Exploring best practices using ITS 2.0 data categories

- Improving real-time translation and multilingual publishing processing by applying extensions, e.g. Readiness:
  - ITS 2.0 extension data category proposal.
  - Linguaserve is applying Readines in  both use cases involved:
    - Applied in CMS-TMS showcase (WP3, poster 3)
    - Applicability in Online Translation system (WP4)
  - It indicates the readiness of a document for submission to L10n processes or provides an estimate of when it will be ready for a particular process.
  - It can be used in expert systems for automatic processing.

# (5) Next steps and some proposals

- Training and methodologies
  - Pre-editing: ITS2.0 usage and training kits.
  - EDI-TA: Post-editing contextual, activation and identification rules.

- Specific tools
  - Pre-editing:
    - Full HTML5 compliance and ITS2.0 annotation facilities
    - Writing tools for content quality, and controlled language for post-editing output adaptation.
  - Post-editing:
    - Specific language-dependent and language-independent post-editing rules and functionalities.
    - ITS 2.0 assistance and viewing functions for post-editors.

**ITS 2.0**

Agencia Tributaria

**MultilingualWeb Workshop**
**Making the Multilingual Web Work**
**Rome, 12–13 March 2013**

**www.agenciatributaria.es**