

CORTEGAL. Corpus de textos galegos escritos por estudantes no ámbito académico

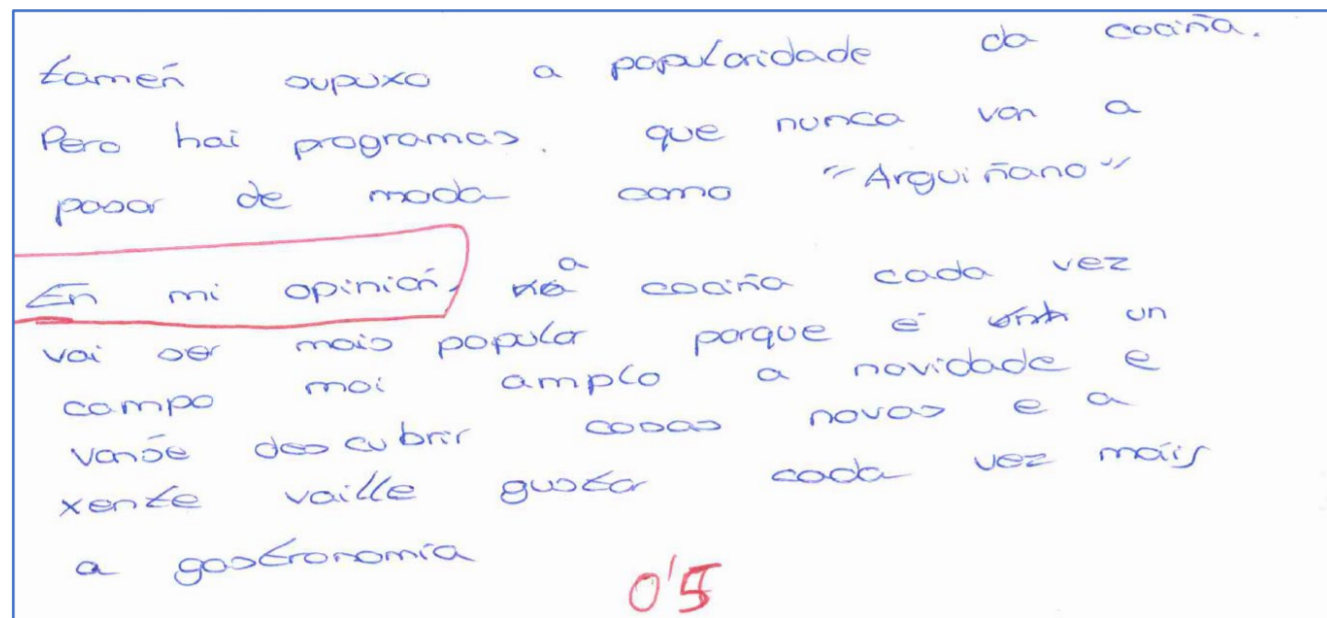
María Álvarez de la Granja
Instituto da Lingua-USC
maria.alvarez.delagranja@usc.gal



I Xeira CLARIAH-GAL
Santiago de Compostela
2 de maio de 2024

1. Presentación

- ✓ 1000 textos argumentativos dos exames de lingua galega das probas ABAU (curso 2016-2017), con anotación das formas non estándares.
- ✓ Obxectivos:
 - Coñecer as características da produción escrita en lingua galega do alumnado de Galicia ao finalizar a educación secundaria.
 - Servir como ferramenta de axuda na aula.



2. Transcripción

- ✓ Transcripción e anotación en TEITOK.
- ✓ Transcripción e etiquetado das formas riscadas polo/a estudante, das claramente engadidas a e das de lectura dubidosa.

```
>Arguiñano/<tok><tok id="w-179" lemma="mais" pos="RG"></tok><tok id="w-219" dcfom="." problem="D_pm_om"></tok></p><p id="p-4"><tok id="w-180" gcfom="Na" problem="G_det_om" psource="G_sp" lemma="en" pos="SP">En/<tok><tok id="w-181" lcfom="miña" olemma="mi" problem="L_w_su" psource="L_sp" lemma="meu" pos="DP1FS">mi/<tok><tok id="w-182" lemma="opinión" pos="NCFS000">opinión/<tok><tok id="w-183" lemma="e" pos="F">e/<tok><tok id="w-184" lemma="e" pos="Fg">e/<tok></del>[add:tok form="a" id="w-185" lemma="o" pos="DA0FS0">a/<tok>[add:tok id="w-186" lemma="cocina" pos="NCFS000">cocina/<tok><tok id="w-187" lemma="cada" pos="D10NN0">cada/<tok><tok id="w-188" lemma="vez" pos="NCFS000">vez/<tok><tok id="w-189" lemma="e" pos="F">e/<tok><tok id="w-190" lemma="ser" pos="VSN000">ser/<tok><tok id="w-191" ocfom="mais" problem="O_ac_om" lemma="ir" pos="VMIP350">vai/<tok><tok id="w-192" lemma="ser" pos="VSN000">ser/<tok><tok id="w-193" lemma="por" problem="O_ac_om" lemma="mais" pos="RG">mais/<tok><tok id="w-194" lemma="popular" pos="AQ0CS">popular/<tok><tok id="w-195" lemma="e" pos="F">e/<tok><tok id="w-196" lemma="un" pos="VSP350">é/<tok><del>[tok form="e" id="w-197" lemma="un" pos="Fg">un/<tok></del>[tok id="w-198" lemma="un" pos="VSP350">un/<tok><tok id="w-199" lemma="campo" id="w-197" lemma="campo" pos="NCMS000">campo/<tok><tok id="w-198" lemma="moi" pos="RG">moi/<tok><tok id="w-199" lemma="e" pos="F">e/</del>
```

3. Anotación

- ✓ Sistema multicapa de seis niveis lingüísticos: ortográfico, morfolóxico (flexión), léxico, gramatical (sintáctico), semántico e discursivo.
- ✓ Códigos que describen o tipo de desviación respecto ao estándar: omisión de acento gráfico (O_ac_om); substitución dunha palabra estándar por unha non estándar (L_w_su)...
- ✓ Nalgúns casos, códigos que identifican a orixe da diverxencia: transferencias léxicas do español (L_sp); flexión dialectal (M_gal)...
- ✓ Anotación e corrección dos tokens a través dun formulario integrado en TEITOK e anotacións multipalabra mediante arquivos *standoff*.
- ✓ Lematización (lema e categoría gramatical) mediante *Freeling* con posterior revisión manual.
- ✓ Asignación ás formas léxicas non estándares (p.e. *platos*) dun lema estándar (*prato*) e dun lema orixinal (*plato*).

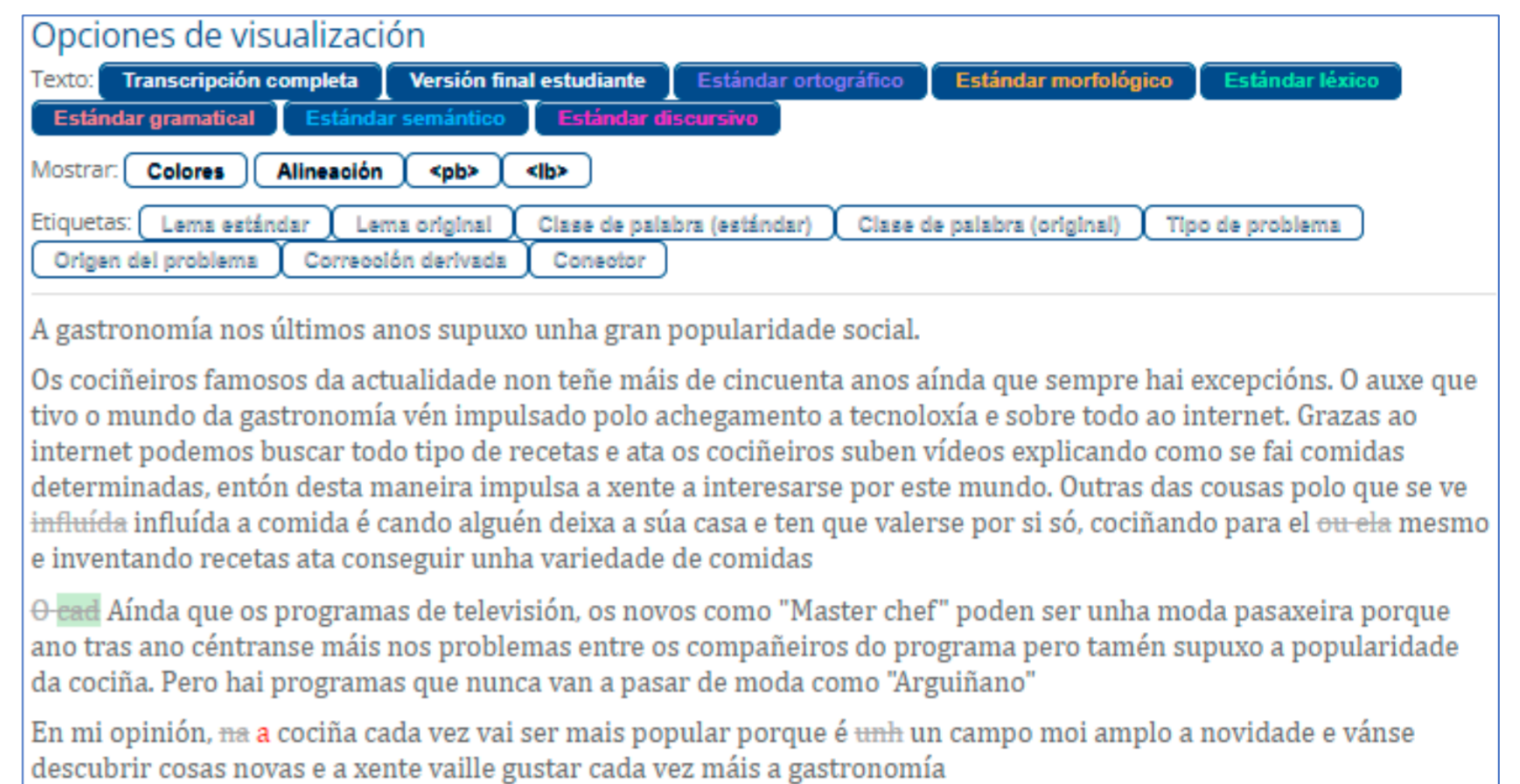
Token value (w-191): mais	
pform	Transcription (inner XML)
form	Student final version
ocform	Orthographic standard
mcfom	Morphological standard
lcfom	Lexical standard
gcfom	Grammatical standard
scfom	Semantic standard
dcfom	Discursive standard

Anotación multipalabra	
lemma	Standard lemma
olemma	Original lemma
pos	POS tag (standard)
opos	POS tag (original)
problem	Type of problem
psource	Source of the problem
dcorrection	Derived correction
arg	Connector

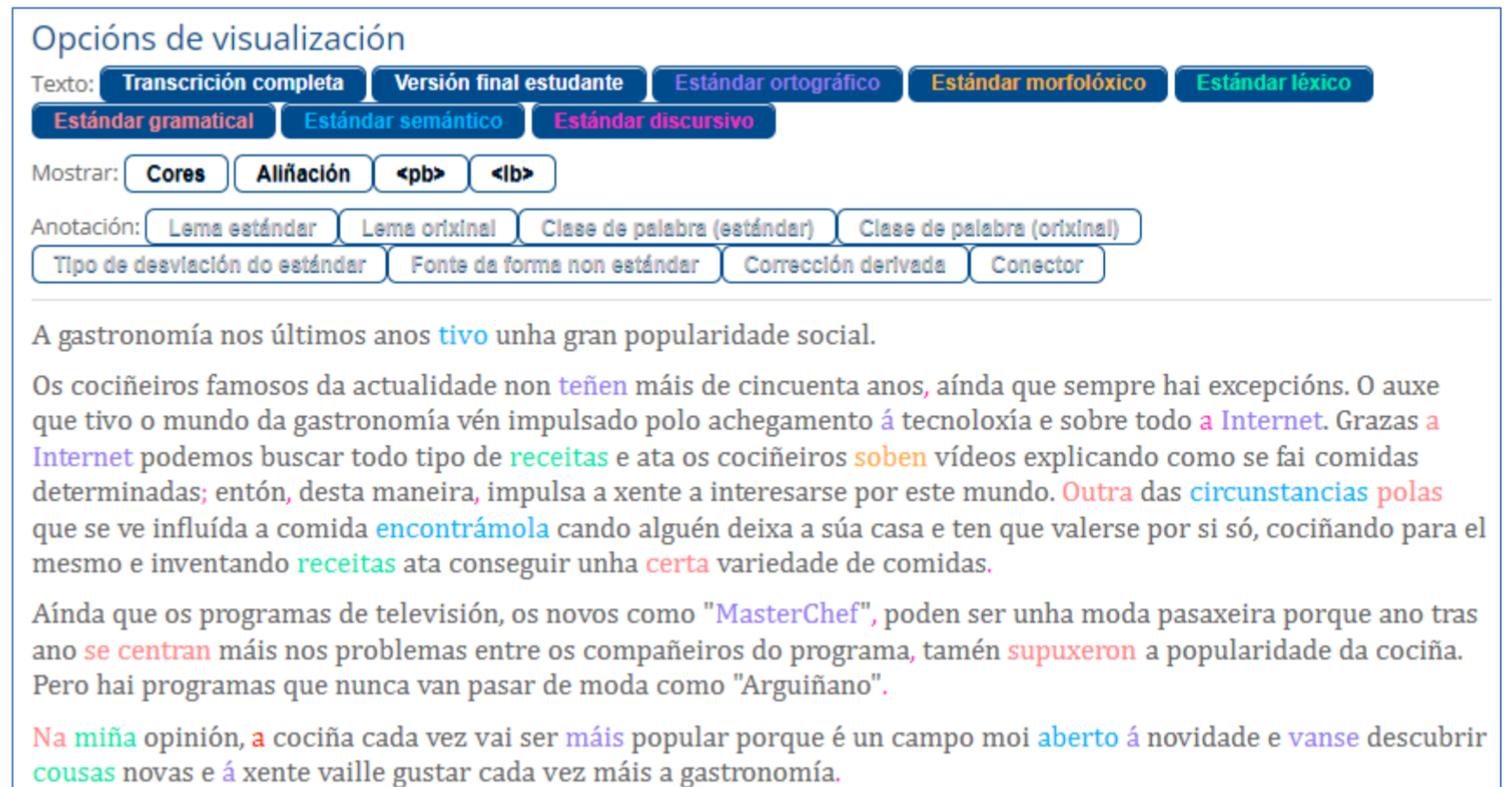
- ✓ Asignación de metadatos cuantitativos: número de lemas, palabras, enunciados e parágrafos, densidade léxica, media de enunciados por parágrafo e de palabras por enunciado e número de palabras do enunciado máis curto e do máis longo.

4. Visualización

- ✓ Distintas posibilidades de visualización:
 - Coas formas eliminadas (en gris e riscadas), como na imaxe, ou na versión final, sen as formas suprimidas. As engadidas destacan en vermello e as de lectura dubidosa con fondo verdoso.



- Coas estandarizacións de cada nivel destacadas en distintas cores: no último nivel, o discursivo, visualízanse todas elas, como na imaxe.



5. Consultas

- ✓ Entre outras, buscas polos códigos das formas desviantes do estándar.
- ✓ Uso dos metadatos como filtros de documentos.

- ✓ Concordancia KWIC con diferentes ordenacións posibles.
- ✓ Consulta por frecuencia, que permite obter datos de distribución do elemento buscado de acordo con diferentes criterios.

Distribución no corpus	
Consulta	Tipo de desviación do estándar = (*,?)L_ac_su(.,*)?
Criterio de agrupamento	Lema orixinal
Total	24
Tamaño de referencia	274617
Gráfica:	Táboa
Contar:	Contar
Descargar:	Seleccionar
élite	7 25.49 29.17
textil	5 18.21 20.83
canón	4 14.57 16.67
periodo	3 10.92 12.5
colón	2 7.28 8.33
Nóbel	1 3.64 4.17
nóbel	1 3.64 4.17
vértice	1 3.64 4.17