

G2P-DLP: An Open-Source System for Grapheme-to-Phoneme Conversion

Leonor Martins^{1,2}, Ana Salgado^{2,1} & Carlos Silva^{2,3}

Institute of Lexicology and Lexicography of the Portuguese Language, Academy of Sciences of Lisbon (ILLP-ACL)¹

Centre of Linguistics of University of Porto (CLUP)²

Wikimedia Portugal³

3rd International Conference on Data & Digital Humanities

27th-28th November 2025

ELACH, University of Minho, Portugal

Agenda

DDHUM25

1. Introduction
2. G2P Systems
3. Methodology Overview
 - 3.1. Data
 - 3.2. Regular Expressions
 - 3.3. Pipeline Architecture
 - 3.4. Open-Source Implementation
4. Results and Discussion
 - 4.1. Level 1
 - 4.2. Levels 2 and 3: Scalability and Sustainability
5. Future Applications and Implications
6. Conclusion
- References

1. Introduction



European Portuguese currently lacks openly accessible grapheme-to-phoneme (G2P) tools, limiting access to phonological data and thus hindering phonological research and education



This work addresses this gap by introducing a G2P system:

Findable

Accessible

Interoperable

Reusable

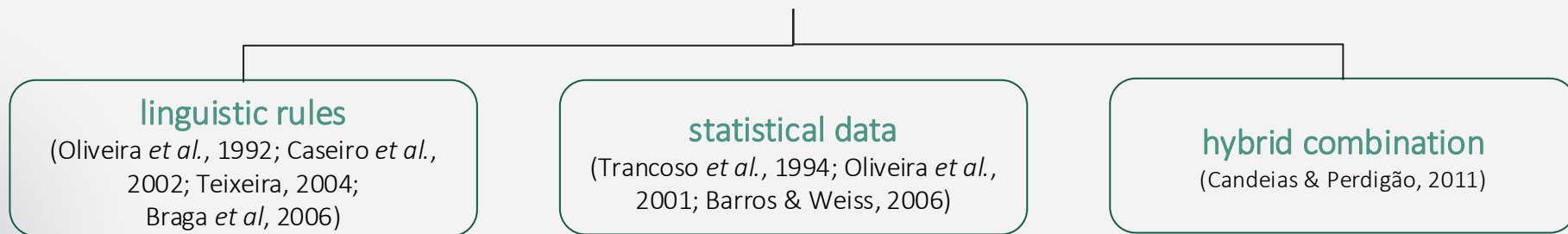
2. G2P Systems

Traditionally, G2P (grapheme-to-phoneme) systems convert written text into phonetic representations.

But also...



Efforts to develop these systems for **European Portuguese** have produced a range of algorithms and models based on...



2. G2P Systems

DDHUM25

Grafone

<https://portulanclarin.net/workbench/it-grafone/>

absence of information about the phonetic conventions and characters (*e.g.*, representation of the central vowel [i] as /ə/)

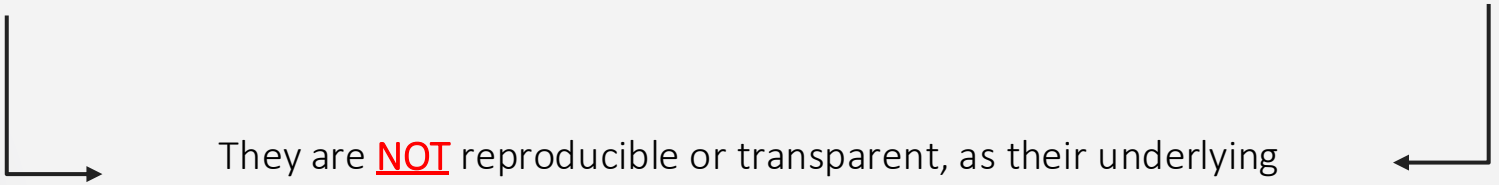
incorrections in the transcription of certain vowels (*e.g.*, bebé, */bəb'ɛ/)

Convert Text to IPA Transcription

<https://european-portuguese.info/ipa#O%3%A1%20mundo!>

inconsistencies on the transcription of certain vowels (*e.g.*, amarelo *[e.me.'re.lu]) and the diphthong [ɔj] (*e.g.*, comboio *[kõ.'boj.u])

incorrections in the transcription of certain vowels (*e.g.*, normal *[nur.'mat])



They are **NOT** reproducible or transparent, as their underlying methodologies and datasets remain undisclosed to the public

3. Methodology Overview

3.1. Data

To develop and evaluate our system, we adopted a **dictionary-based approach**, which relies on a list of words (a **lexicon**) corresponding to the **lemmas** of the *Dicionário da Língua Portuguesa* of the Academy of Sciences of Lisbon (DLP-ACL)

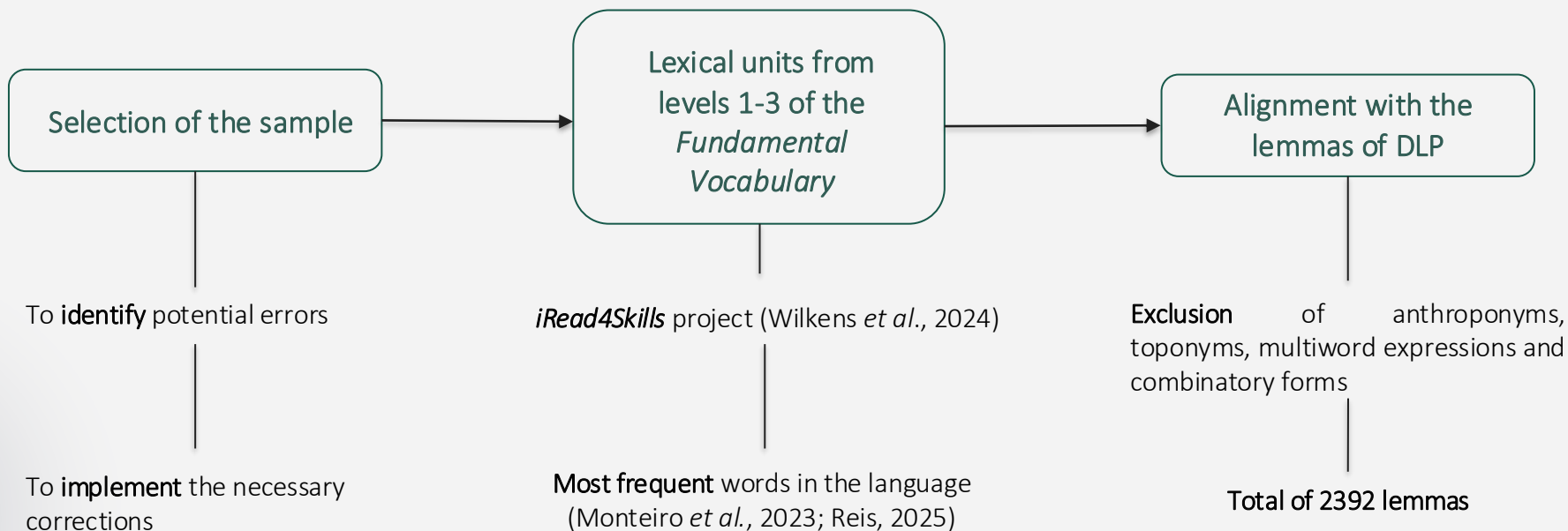
	A	B	C	D
1	aberto			
2	alegre			
3	anterior			
4	atento			
5	barato			
6	básico			
7	brilhante			
8	central			
9	certo			
10	clássico			
11	colorido			
12	completo			



<https://dicionario.acad-ciencias.pt/>

3. Methodology Overview

3.1. Data



3.2. Regular Expressions

Regular expressions (regex) are versatile search patterns, effective in identifying and/or modifying data based on predefined character sequences
(Goyvaerts, 2007; Goyvaerts & Levithan, 2009)

broad applicability

high interoperability

great stability across software environments

Sample of the syntax of regex characters

Character	Description	Example
^	Selects the initial position	^g → ganhar
\$	Selects the final position	á\$ → cá
()	Selects the character(s) on the left OR the character(s) on the right	(m n lh) → malha; mana; mala
(?=) (positive lookahead)	Selects the character when it is followed by what goes after =	p(?=ão) → pãõ; pã
(?!) (negative lookahead)	Selects the character when it is NOT followed by what goes after !	p(?!e) → perto; preto; porta
(?<=) (positive lookbehind)	Selects the character when it is preceded by what goes after =	(?<=o)l → voltar; possível
(?<!) (negative lookbehind)	Selects the character when it is NOT preceded by what goes after !	(?<!g)ui → seguir; aquilo

3.2. Pipeline architecture

DDHUM25

```
(?<=n)(?=[bcçdfgjlmnpqrtvxz]) Aa ab .*
```

AB

- 1 lindo -> lin.do
- 2 longo -> lon.go
- 3 antigo -> an.tigo
- 4 branco -> bran.co
- 5 contente -> con.ten.te
- 6 diferente -> diferen.te

```
(?=^[áâêéèêñíóòôúùû]) Aa ab .*
```

AB

- 1 ótimo -> 'ótimo
- 2 último -> 'último
- 3 água -> 'água
- 4 área -> 'área
- 5 árvore -> 'árvore
- 6 único -> 'único

```
^j Aa ab .*
```

AB

3

- 1 ja'neiro -> ʒa'neiro
- 2 jar'dim -> ʒar'dim
- 3 ja'nela -> ʒa'nela
- 4 joga'dor -> ʒoga'dor
- 5 jo'gar -> ʒo'gar
- 6 jun'tar -> ʒun'tar

Syllabic Division

- selects a **boundary** between syllables
- inserts a **syllabic divider**

Stress Marking

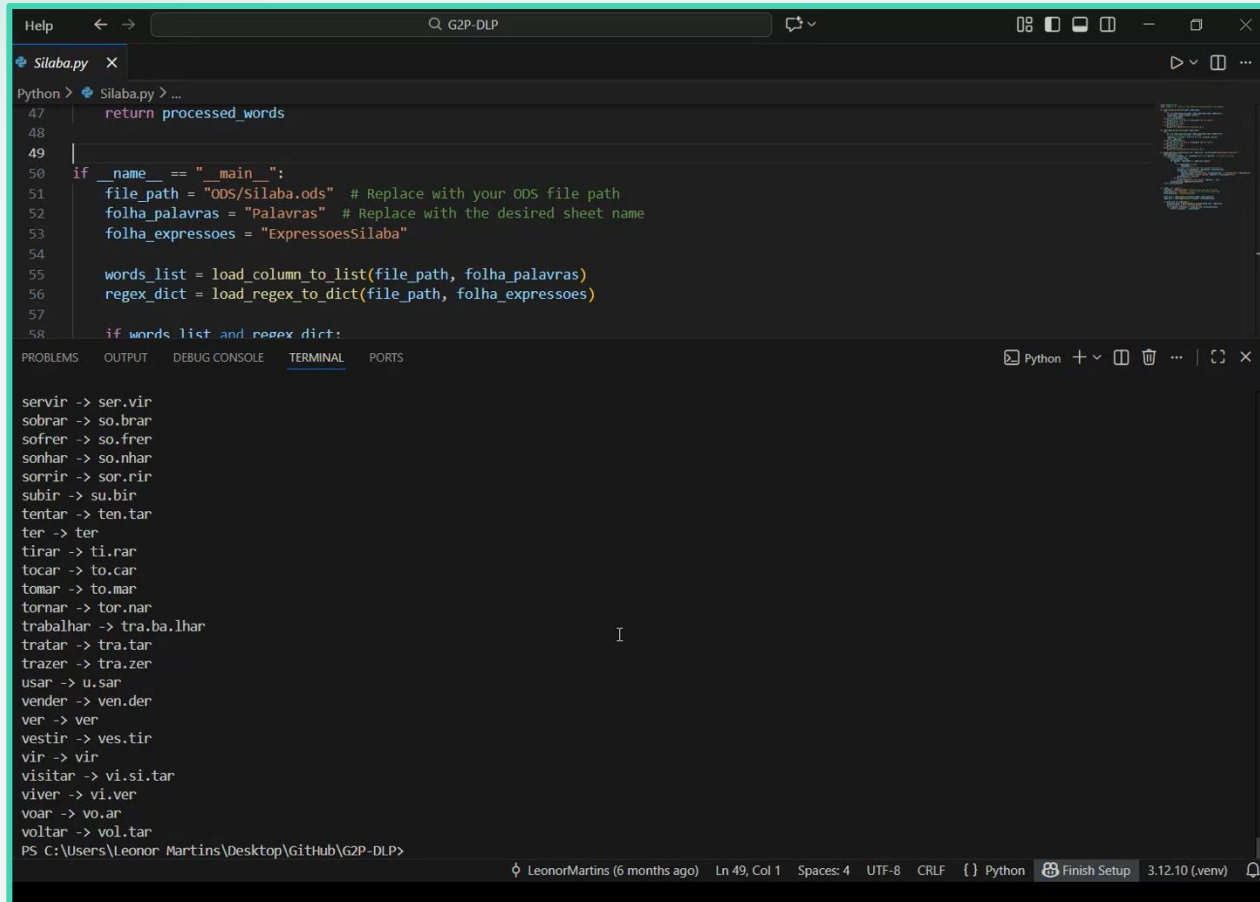
- selects a **position** before a syllable or a **boundary** between syllables
- inserts a **stress marker**

Phonetic Transcription

- **replaces** a segment with another segment

3.2. Pipeline architecture

DDHUM25



```
Help  ← →  G2P-DLP  [Icons]  [Close]  [Maximize]  [Fullscreen]  [Refresh]  [Close]
```

```
Silaba.py  [Close]
```

```
Python  >  Silaba.py  >  ...
```

```
47     return processed_words
48
49
50 if __name__ == "__main__":
51     file_path = "ODS/Silaba.ods" # Replace with your ODS file path
52     folha_palavras = "Palavras" # Replace with the desired sheet name
53     folha_expressoes = "ExpressoesSilaba"
54
55     words_list = load_column_to_list(file_path, folha_palavras)
56     regex_dict = load_regex_to_dict(file_path, folha_expressoes)
57
58     if words_list and regex_dict:
```

```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL  PORTS  [Python]  [+/-]  [Icons]  [Close]
```

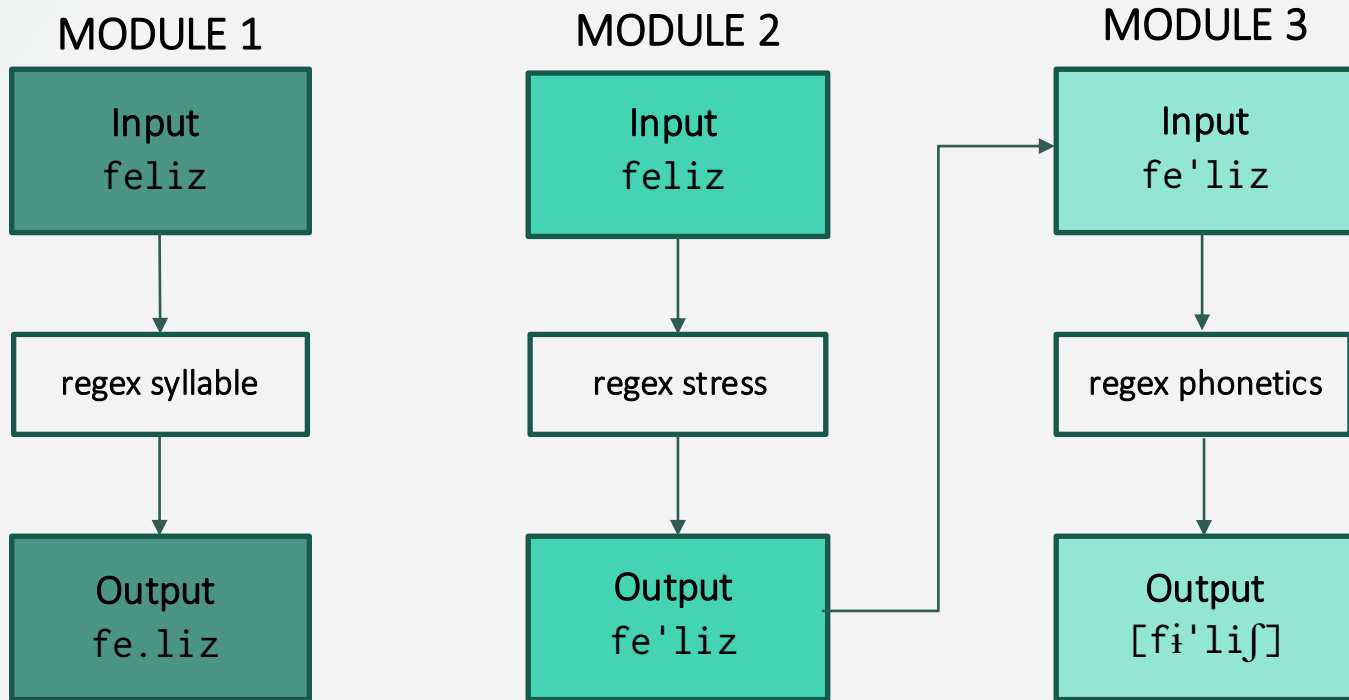
```
servir -> ser.vir
sobrar -> so.brar
sofrer -> so.frer
sonhar -> so.nhar
sorrir -> sor.rir
subir -> su.bir
tentar -> ten.tar
ter -> ter
tirar -> ti.rar
tocar -> to.car
tomar -> to.mar
tornar -> tor.nar
trabalhar -> tra.ba.lhar
tratar -> tra.tar
trazer -> tra.zer
usar -> u.sar
vender -> ven.der
ver -> ver
vestir -> ves.tir
vir -> vir
visitar -> vi.si.tar
viver -> vi.ver
voar -> vo.ar
voltar -> vol.tar
PS C:\Users\Leonor Martins\Desktop\GitHub\G2P-DLP>
```

```
LeonorMartins (6 months ago)  Ln 49, Col 1  Spaces: 4  UTF-8  CRLF  [Python]  [Finish Setup]  3.12.10 (venv)
```

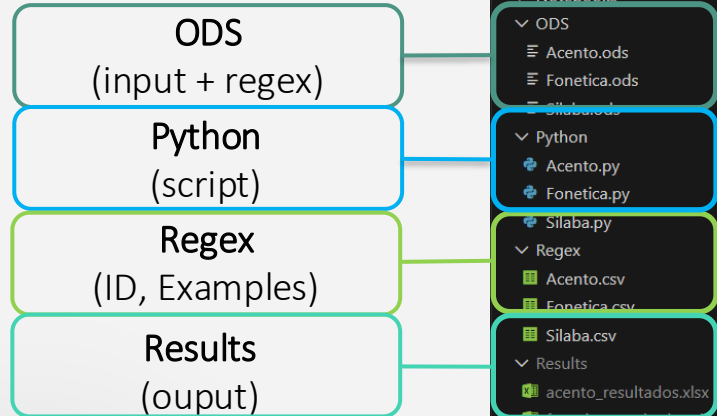
```
servir -> ser.vir
sobrar -> so.brar
sofrer -> so.frer
sonhar -> so.nhar
sorrir -> sor.rir
subir -> su.bir
tentar -> ten.tar
ter -> ter
tirar -> ti.rar
tocar -> to.car
tomar -> to.mar
tornar -> tor.nar
trabalhar -> tra.ba.lhar
tratar -> tra.tar
trazer -> tra.zer
usar -> u.sar
vender -> ven.der
ver -> ver
vestir -> ves.tir
vir -> vir
visitar -> vi.si.tar
viver -> vi.ver
voar -> vo.ar
voltar -> vol.tar
```

3.2. Pipeline architecture

DDHUM25



3.4. Open-Source Implementation



```
Python > Silaba.py > load_column_to_list
1 import re
2 import unicodedata
3 from pathlib import Path
4 import pandas as pd
5 import openpyxl # writer engine
6
7 # Do not alter the first chunk if it matches one of these (dots removed, case-insensitive)
8 PROTECTED_FIRST_CHUNKS = {"afro", "alto", "amarelo", "amor", "anglo", "ãntero", "ásio", "austro",
9
10 def load_column_to_list(file_path, sheet_name):
11     """Read first column from a sheet into a list of strings."""
12     try:
13         df = pd.read_excel(file_path, sheet_name=sheet_name, header=None)
14         return df[0].dropna().astype(str).tolist()
15     except FileNotFoundError:
16         print(f"Error: The file at {file_path} was not found.")
17     except ValueError as e:
18         print(f"Error: {e}")
19     except Exception as e:
20         print(f"An unexpected error occurred: {e}")
21
22 def load_regex_to_dict(file_path, sheet_name):
23     """Read two columns into a dict {pattern: replacement} as strings."""
24     try:
```

3.4. Open-Source Implementation



Repository Structure

- Debug/
 - └─ Log and debug files generated during the execution of the scripts and notebooks.
- Notebooks/
 - └─ Jupyter notebooks corresponding to the three modules (syllable, stress, phonetics). They present the process of developing the Python files and the results obtained.
- Python/
 - └─ Standalone Python scripts for the three main modules.
 - silaba.py
 - acento.py
 - fonetica.py
- ODS/
 - └─ ODS-format data files containing the corpora used in the modules. These files specify the contexts, patterns, and substitutions applied.
- Regex/
 - └─ CSV tables with the regular expressions used in each module. Each line includes:
 - Regex ID
 - Context Pattern
 - Regex expression (Find)
 - Applied substitution (Replace)
 - Example of application

4. Results and Discussion

DDHUM25

4.1. Level 1

Syllabic Division

100%

All syllable boundaries were detected:

- C.C (partir → par.tir)
- V.C (tratar → tra.tar)
- G.C (baixo → bai.xo)
- G.V (meia → mei.a)
- V.V (voar → vo.ar)

Stress Marking

99.6%

All stress patterns were detected:

- **proparoxytones** (férias → 'férias)
- **paroxytones** (festa → 'festa)
- **oxytones** (flor → 'flor)
Except *simples*, *apenas* and *menos*:
- *sim'ples → 'simples
- *ape'nas → a'penas
- *me'nos → 'menos

Phonetic Transcription

97.6%

All consonantal segments were correctly

transcribed ([pu'xar] → [pu'ʃar]).

The vowel segments were more problematic:

- *bebé*: *[bɨ'bɛ] → [bɛ'bɛ]
- *normal*: *[nur'maɫ] → [nɔr'maɫ]
- *ação*: *[e'sẽw̃] → [a'sẽw̃]
- *amanhã*: *[eme'ɲẽ] → [ame'ɲẽ]

4.2. Levels 2 and 3: Scalability and Sustainability

DDHUM25

Syllabic Division

99.6%

A small set of **ambiguous graphemic sequences** challenged syllabification, namely <ai> and <ui>:

- rainha (ra.i.nha) vs. baixo (bai.xo)
- fluidez (flu.i.dez) vs. cuidar (cui.dar)

Stress Marking

99.5%

The few deviations were limited to **predictable sources of ambiguity**:

- ambiguous graphemic sequences (*'juiz → ju'iz)
- double marked words (*'ór'fão → 'órfão)
- exceptions (*am'bos → 'ambos)

Phonetic Transcription

97.7%

Errors were restricted to a small number of **well-defined contexts**:

- grapheme <x> (*exame*: *[e'jami] → [e'zami])
- adverbs in *-mente* (*somente*: *[su'měti] → [so'měti])
- exceptions of vowel reduction (*coleção*: *[kulɨ'sěw] → [kulɛ'sěw])

5. Future Applications and Implications

DDHUM25



Technical Framework and Sustainability

- uses **open-source, widely supported technologies**
- provides a **transparent, replicable, and extensible** framework
- flexible architecture that supports **easy updates** and adaptation to **new corpora or languages**



Pedagogical Applications

- can be integrated into tools for **pronunciation training** and **language learning**
- supports **literacy development, first and second-language teaching** and the **creation of educational materials**



Educational Resource Development

- enables rapid production of **annotated datasets** that can be used in classroom resources
- supports teachers and learners through **standardized phonological representations**



Research and Linguistic Innovation

- facilitates **large-scale phonological analysis**
- could support **comparative studies on prosody** and serve as a preprocessing component for **speech technologies**
- contributes to **language conservation and revitalization**

6. Conclusion

Our G2P system demonstrates **stability** and **accuracy** across various modules and dataset sizes



Providing **openly available** solutions is crucial to enabling **equitable and democratic participation in linguistic research and development**



This project offers a **FAIR resource** that illustrates the **impact of open infrastructures**

The screenshot shows the website interface for the Dicionário da Língua Portuguesa. At the top left is the logo of the Academia das Ciências de Lisboa. At the top right is the official seal of the Academia das Ciências de Lisboa. The main heading is 'igualdade'. Below it, there is a search bar containing 'i.gual.da.de' and a pronunciation guide icon with the phonetic transcription [igˈwaˈɫɔdɐdɨ]. To the right of the search bar is a small red box with the text 'DLPC 2001'. Below the heading, the word is identified as 'nome feminino'. The first definition is '1. qualidade ou estado do que não apresenta diferença; relação de paridade entre coisas ou seres iguais'. Below this, there is a list of synonyms: 'igualdade de circunstâncias; igualdade de direitos; igualdade de oportunidades'. At the bottom, there are sections for 'SINÓNIMOS' (identidade) and 'ANTÓNIMOS' (desigualdade; diferença).

References

- Aho, A. V. (1991). Algorithms for finding patterns in strings, *Handbook of theoretical computer science (vol. A): algorithms and complexity*. MIT Press, Cambridge, MA.
- Braga, D., Coelho L., Resende Jr, F. G. (2006). A Rule-Based Grapheme-to-Phone Converter for TTS Systems in European Portuguese. Em *VI International Telecommunications Symposium*, 328-333.
- Candeias, S. & Perdigão, F. (2011). Integração linguística em sistemas de conversão de grafema para fone[ma]. Em A. R. Luís (Ed.), *Estudos de linguística*, 1, 181-194. doi: 10.14195/978-989-26-0231-8_12
- Candeias, S. & Perdigão, F. (2013). Grafone: uma ferramenta de/com recurso à informação linguística. Em *Textos Seleccionados, XXVIII Encontro Nacional da Associação Portuguesa de Linguística*, 189-203. doi: <http://doi.org/10.1109/its.2006.4433293>
- Caseiro, D. A., Trancoso, I., Oliveira, L. & Viana, C. (2002). Grapheme-to-phone using finite-state transducers. Em *Proceedings of the IEEE Workshop on Speech Synthesis*, 215-218. doi: 10.1109/WSS.2002.1224412
- DLP-ACL. (2025). Academia das Ciências de Lisboa. *Dicionário da Língua Portuguesa*. Salgado, A. (Coord.). Lisboa: Academia das Ciências de Lisboa. Disponível em <https://dicionario.acad-ciencias.pt/>
- Goyvaerts, J. (2007). *Regular Expressions: The Complete Tutorial*. LuluPress, Incorporated.
- Goyvaerts, J. & Levithan, S. (2009). *Regular Expressions Cookbook*. O'Reilly Media, Inc.
- IPA – International Phonetic Association. (n.d.). *International Phonetic Alphabet*. Retirado de <https://www.internationalphoneticassociation.org/>.
- Martins, L. (2025). *Desenvolvimento de um sistema de conversão grafema-fone[ma] para o vocabulário Fundamental da Academia das Ciências de Lisboa*. Dissertação de Mestrado, Faculdade de Letras da Universidade do Porto. <https://repositorio-aberto.up.pt/handle/10216/168957>
- Martins, L. (2025). *G2P-DLP*. <https://github.com/LeonorMartins/G2P-DLP>
- McKinney, W. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference*, Vol. 445, pp. 51-56.
- Monteiro, R., Amaro, R., Correia, S., Pintard, A., Gauchola, R., Moutinho, M., & Blanco Escoda, X. 2023. *iRead4Skills – Complexity Levels*. Zenodo. <https://doi.org/10.5281/zenodo.10459090>
- Oliveira, L. C., Viana, M. & Trancoso, I. (1992). A rule-based text-to-speech system for portuguese. Em *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2, 72-76. 10.1109/ICASSP.1992.226117 Imprensa da Universidade de Coimbra. doi: [10.1109/ICASSP.1992.226117](https://doi.org/10.1109/ICASSP.1992.226117)
- Reis, M. L. (2025). *Vocabulário Fundamental da Academia das Ciências de Lisboa: seleção lexical, alinhamento de sentidos e codificação*. Tese de Mestrado. Faculdade de Ciências Sociais e Humanas da Universidade NOVA de Lisboa.
- Teixeira, J. (2004). *A Prosody Model to TTS Systems*. [Tese de Doutoramento, Faculdade de Engenharia da Universidade do Porto]. <http://hdl.handle.net/10198/1496>
- Trancoso, I., Viana, M. C., Silva, F. M., Marques, G. C. & Oliveira, L. C. (1994). Rule based vs neural network-based approaches to letter-to-phone conversion for Portuguese common and proper names. Em *3rd International Conference on Spoken Language Processing (ICSLP 1994)*, 1767-1770. doi: 10.21437/ICSLP.1994-197
- Trigo, L. & Silva, C. (2022). Comparing lexical and usage frequencies of palatal segments in Portuguese. In *Computational Processing of the Portuguese Language: 15th International Conference, PROPOR 2022, Fortaleza, Brazil, March 21–23, Proceedings*, pages 353–362, Berlin. Springer.
- Wilkens, R., Pintard, A., François, T., Barbosa, S., Reis, M. L., Amaro, R., Ribeiro, E., Mamede, N., Baptista, J., Blanco, X., Catena, A., Gauchola, R., & Mu, K. (2024). *iRead4Skills – Basic Lexicons per Complexity Level (v1.0) [Data set]*. Zenodo. <https://doi.org/10-81/zenodo.10889986>

G2P-DLP: An Open-Source System for Grapheme-to-Phoneme Conversion

Leonor Martins^{1,2}, Ana Salgado^{3,4} & Carlos Silva^{2,5}

Institute of Lexicology and Lexicography of the Portuguese Language, Academy of Sciences of Lisbon (ILLP-ACL)¹

Centre of Linguistics of University of Porto (CLUP)²

Academy of Sciences of Lisbon (ACL)³ & Faculty of Arts and Humanities of University of Porto (FLUP)⁴

Wikimedia Portugal⁵

3rd International Conference on Data & Digital Humanities

27-28 November 2025

ELACH, University of Minho, Portugal