# Rule-based Surface Realization of Roman...

Ciprian-Virgil Gerstenberger
UiT The Arctic University of Norway
ciprian.gerstenberger@gmail.com

## Abstract

Due to its reliance on context and intricate grammatical rules, the Romanian weak pronoun system presents a challenge not only for language learners, both native and non-native speakers, but also for linguistic description and computational processing.

The present work addresses the challenges of Romanian weak pronouns from a computational processing perspective. Accordingly, it has three main goals: (1) to present the implementation of a rule-based model for generating contextually accurate surface forms of Romanian weak pronouns, (2) to describe the compilation of a database of relevant and contrasting surface realizations, and (3) to test the effectiveness of the model.

This serves as a proof of concept, demonstrating both the transparency and the effectiveness of the model based on an appropriate linguistic description.

## 1 Introduction

Romanian weak pronouns (henceforth, RWPs), commonly referred to as RWPs, exhibit a variety of surface forms governed by intricate, contextually dependent morpho-phonological rules. As RWPs, they pose challenges not only for linguistic description but also for computational natural language processing (henceforth, NLP).

Despite extensive research on Romance clitics within various theoretical frameworks – including Generative Grammar, as explored by Dobrovie-Sorin (1979), Theophanopoulou ...

(2007), and Săvescu Ciucivara ... cal Functional Grammar (LFG), ... by Bárbu and Toivonen (2018); the Romanian ... Phrase Structure Grammar (HPS... cussed in Monachesi (2001) and ... (2005); various Optimality The... dels advanced by Popescu (200... Ciulianu (2000), Legendre (2... (2003), and (2014); and Dynami... Syntax, as promoted by Klein (...

... demanded for practical, ... processing that can accurately gene... their usage in real-world cont...

The current study addresses t... cusing on the creation of a rul... designed to handle the surfac... RWPs, tested using a specializ... base. This, in turn, contribu... inputs for enhancing computat... applications for Romanian.

## 2 Romanian Weak Pronouns

As with both other Romance language... morphological system t... both strong pronouns and weak (c... nouns. Romanian has a complex ... with fully marked case, number... distinction, comprising both ... referred to as RWPs. Moreover, there ... form and a form gover(henceforth, ... of the preceding morpho-phonological ... As RWPs, they appear in both s... additional forms. The first syllab... overall morphological requirement for as... categorys (and hence) for the possibil... (optional sandhi), is determi... each Romance. Sandhi r... logical adjustments that occu... word or across... influenced... the syllable (and thus...

forms are referred to as sandhi forms.

RWPs occur in a fixed order [RWP_clit/cus < RWP_dat < RWP_acc] in clitic clusters (or sequences), alongside other clitics such as negation, auxiliary verbs, or adverbial particles (ex. 1). The cluster can occur both in preverbal position (e.g., in declarative sentences, ex. 3), with the and adverbial particles which can occur only preverbally.

(1) Nu ni le- ai
not cl_1.pl cl_3.pl.acc have_2.sg.pres
mai dat.
more given
«You didn't give them to us anymore.»

(2) Mi le dai
cl_1.sg cl_3.pl.acc give_2.sg.pres
acum.
now
«You give them to me now.»

(3) Dă- mi -le acum!
give_2.sg.imp cl_1.sg cl_3.pl.acc now
«Give them to me now!»

(4) Le dai mere.
cl_3.pl.acc give_2.sg.pres apples
«You give apples to them.»

Romanian exhibits various types of ambiguity that complicate the interpretation of the RWP data, such as case syncretism between accusative (ex. 2) and dative (ex. 4) plural, part-of-speech homonymy (between the dative-reflexive RWP function «and» in ex. 5), phonetic-graphemic ambiguity (syllabic form [mi] in ex. 6 vs. syllabic form with nasal [i] in ex. 7), as well as hyphen ambiguity (marking prosodic clitic-hood, means that they may reduce (asyllabic mi -i in ex. 7 vs. marking post-verbal (syllabic mi -i in ex. 10). Optional sandhi mi -i in ex. 3).

(5) i le cumpără
cl_3.sg.dat cl_3.pl.acc buy_3.sg.pres
și i le
and_conj cl_3.sg.dat cl_3.pl.acc
revinde.
resell_3.sg.pres
«He/she buys them for him-/herself and resells them for him-/herself»

(6) Mi -l dai.
cl_1.sg.dat cl_3.sg.acc give_2.sg.pres
«You give it to me.»

(7) Mi -o dai.
cl_1.sg.dat cl_3.sg.acc give_2.sg.pres
«You give her/it to me.»

(8) ... îmi cartea.
give_2.sg.pres cl_1.sg.dat book_def.sg.acc
«You give me the book.»

(9) Dă- -mi- cartea!
give_2.sg.imp cl_1.sg.dat book_def.sg.acc
«...»

As a general observation, syllabic whether obligatory or optional in the rightmost RWP. Obligatory the following item, to the right by the occurrence of an auxiliary with a vowel (e.g., a) or the 3.sg. RWP -o (form mi -o ex. 7).

Obligatory sandhi to the preceding the left, occurs if the context sandhi to the right is not present form (see Section 3) item in surface. If the preceding item is also an RWP it serves as syllabic rightmost RWP (ex. 6). If the rightmost RWP is the only item in the surface a prosthetic form in preverbal form (î mi in ex. 8). In postverbal position verb functions as the syllabic ex. 9).

The RWP underlying, giving forms exhibit a special behavior as single RWPs, thus in both the and leftmost position, they su- respective parts inter- pretation of the respectively an instance of the aforementioned syncretism (ex. 4). In other clusters, the surface forms re- homonymy (between the underlying forms).

When obligatory sandhi is not specific contextual conditions, in an RWP RWPs can undergo optional means that they may reduce (asyllabic clitic-hood, but such a reduction is not (syllabic ex. 10). Optional sandhi following item, to the right, following item starts with a (ex. 11). Similarly, some RWPs of length one can optionally at- ceding item, to the left, if the ends with an unstressed vowel -/herself and resells (ex. 12). In the same context, the surface as an î -i phonetic form (ex. 13)

(10) Le aduci mere.
cl_3.pl.dat bring_2.sg.pres apples
«You bring them apples.»

(11) Le- aduci mere.
cl_3.pl.dat bring_2.sg.pres apples

«You bring them apples.»

(12) y-reau    s -mi
     want_1.sg ch.acl r_e1s. sg. dat
     dai      mere.
     give_2.sg pl pers es
     «I want you to give me apples.»

(13) y-reau    s îmi
     want_1.sg ch.acl r_e1s. sg. dat
     dai      mere.
     give_2.sg pl pers es
     «I want you to give me apples.»

## 3 Model description

The computational implementation of a theoretical model serves as a bridge between abstract concepts and practical applications, providing a platform to explore, validate, and extend the capabilities of the model. In linguistics, where numerous frameworks construct with diverse theoretical... ting these abstract models into computational terms becomes essential (Bender and Langendoen, 2010)

Although numerous computational tools for linguistics are available – such as Hayes et al. (2013) for OT, Kaplan et al. (2004) for LFG, Copestake (2001) for HPSG, to name a few – none of the theoretical approaches mentioned in Section 1 have attempted to demonstrate how these frameworks can be computationally implemented and systematically validated.

In Gerstenberger (2022) I provided comprehensive description of RWPs and their contexts of occurrence using subjectively testable linguistic features based on their (a)syllabicity. Granted the syllable... is notoriously difficult to define in a universally agreed-upon way – is the a unit of speech production (articulation), perception, or both? (cf. for instance, Cole et al. 2008) but it is widely recognized as a building block in phonological theory...

Unlike Dobrovie-Sorin who Dobrovie-Sorin (2013: p. 266) claim that clitic forms are underlyingly asyllabic, Popescu (2003: p. 154) claim of underlying moraic, or Klein (2007: p. 77) use of cluster-prosthetic forms as input, I propose a model in which all underlying representations are uniformly syllabic (see Table 1).

The model I propose for handling... similar to the generative approach... sky and Halle (1968) and the two... phology in Koskenniemi (1983)... two different levels of representation... a set of rules for mapping between... lying and surface levels.

While the underlying representations... theoretical linguistic entities... are their computational linguistic... parts – concrete strings that...

Given the case syncretism between... tive and dative plural RWP form... in Section 2 (-syl ⇒ c) and the... (syl ⇒ abc a syllabic)... as well a... observation that the obligatory... right has higher precedence than... tory sandhi to the left, a constru... as presented for XFST in Beesle... Beesley (2003) seems better suited... these rules (see for instance,

The key idea when modeling RWP... mena is to account for the positi... in the clitic sequence: is the... most (position or not? Next, it... identify the correct contextual... tinguishing between obligatory... sandhi: the asyllabic RWP form... by vowel loss, hyphen addition... the postverbal position is cha... addition of hyphens between a... clitic sequence that follow th... The following XFST grammar... provides definitions and rules... obligatory sandhi to the right... linguistic features based... that are followed by an auxilia... with a vowel or by the RWP form...

```
define MaRWPI "o" a unit of
define ReduceHighV_RWP "mi" | "i"
  ne OtherHigh_RWP bu tu
define SubstituteHighV_RWP "ni" | "v
define HighV_RWP ReduceHighV_RWP |
  DeleteHighV_RWP | SubstituteHigh
define DeleteLowV_RWP "m " | "se"
define LowV_RWP DeleteLowV_RWP |
define SyllabicC_RWP PerHigh_V_RWP | LowV_
define unspecified_Aux "am" | "ai" |
define Optional_Aux "oi" | "vei"
  Vom
define Initial_Aux | C_Initia
  V_Initia
define Host Host V_Initia
  V_RWP
define RWP V_RWP | Syllabic_RWP;
define Rightmost_RWP Syllabic_RWP
```

| Number | Accusative | | | | | Dative af | | | | Ve |
|---|---|---|---|---|---|---|---|---|---|---|
| | 1p | 2p | 3p.m | 3p.f | 3p.ref | 1p | 2p | 3p.m | 3p.f | 3p |
| Sg | /m/ | /te/ | /lu/ | /o/ | /se/ | /mi/ | /ts/ /su/ | /i/ | /i/ | /i/ |
| Pl | /ne/ | /v/ | /i/ | /le/ | /se/ | /nii/ | /vi-/ | /su/ | /l |

Tabela 1: Underlying forms as input for the surface fo... «to be»

```
    Syllabic_RWP .#. Syllabic_RWP;
define Remove_Vowel
    (DeleteLowV_RWP | DeleteHighV_RWP) -> ("")
    / [aeiou] _ )
    / .#. Is_Obligatory_Host;
define Substitute_Form
    "ni" -> "ne" / _ Rightmost_RWP .#. ||
    "vi" -> "v " / _ Rightmost_RWP .#.
    "li" -> "le" / _ Rightmost_RWP .#. ;
define Add_Right_Hyphen Rightmost_RWP -> ...
    _ Is_Obligatory_Host;
define Asyllabic_Transformation Remove_Vowel
    || Substitute_Form || Add_Right_Hyphen ;
regex Asyllabic_Transformation;
```

- Python has excellent debugging tools.
- Python has a rich ecosystem f... as spaCy, NLTK, Stanza.
- Python syntax is cleaner and ... read... than XFST syntax.
- Python scripts can be easily ... or deployed, or packag...

The task of the surface form g...

First, groups of items which behave the same... are reduced... within its specific context...

[The remainder of this page consists of two heavily overlapping text columns that cannot be reliably separated.]

## 4 Model implementation

While the XFST grammar from the previous section is a pr... istic phenomena described in... translated into XFST rules... plementation has been done... are some non-negligeable...

- Python allows for more expressive rule implementations.

and customized for specific types of phenomena or replaced by the annotation. When run on the input, all annotations and diagnostics information can be queried and checked for consistency, it acts as in RWPs 2 values, allowing operations to be performed and merged because of each token. both for a syllabic and an asyllabic

The input sentence (see A segment in Figure 2) is analyzed using a spaCy pipeline and hi refers to phonologic created specifically for this setup to produce the boundaries shown in the code fragment in Figure 1, features such as Python rules operate on each token in the Doc object), social spaCy Doc object. lects (sociolect, regiolect)

For instance, the code checks whether speech registers current token pairs and compares output targets, these are based fned in Gerstenberger (2007), it on a sign and hi, that can occur in both preverbal and postverbal positions. If this condition is met, a further database entry checks whether the token is an RWP such as in tng, the or an RWV to determine whether it can list what nemuoperations should be performed (1 mso wnu sthe museum!»): bo string. Along with default value checks, a form of the current token – such as in its surface that con tion and contextual elements are customizing options sa tions token. _., token_p_., is en WV access i i da («Where are you anis_rightmost_in_clause giving (token m) her apples?») in been implemented for this task. To address hsheots her the dative/accusative case syncretism in the only plural in term forming has been one of the different targets. employed. The entire rule set is freely available comple ble from https://github.com/ciprian-NO/ic positioning, rb_rwp_generation/blob/main/generate_heir syntactic behavior, and interactions with _rwp_surface_forms.ipy in the sentence. When building

To test the effectiveness of the module, we test input database for testi created a test input database covering all re alization of RWPs, several key levant RWP phenomena described in Section structural factors (see Secti 2. sidered, such as: person, numb case of RWPs, cluster length,

## 5 Test database creation
The test input database has b Evidently, the design of the test database on examples fro base is driven by the RWP surf (2022) and the relevant l model proposed. The structure of RWP examples that we entry is a tuple <INPUT, TARGE..., plete sentences have been mini TARGET>, where the input string contains sentence (in a da tains the syllabic input a i.e. entry) underlying by a representation of RWPs (e.g. the «beauty con - Ceuomp mers - l tains all possible output for example Bayits quick lam). Some R ple, for the cursîenliken expres se repeated in different s-ar toate(«May all his/her phone realization in cases o stifen!»), then input is Teh database contains 352 inp toate oa,scrlea]ted by remo470ntgarighthys-entences, distri phens that indicate post25ersbati twi toh rnsey tlarget outpu city. Furthermore, the to boltiagra geotry aaria ylnltab,i7csets wit

```python
# 1. current token is a Linear Order Part (LOP), i.e., can occur both pre- and post-verbally
if token._.is_linear_order_part:
    is_postverbal_token = not is_preverbal(token) and not r_nbor.pos_ == 'VERB'
    if is_postverbal_token:
        postverbal_hyphen = HYPHEN_TOKEN.text

    # 1.1 current token is a Romanian Weak Pronoun (RWP) or a Romanian Weak Verb (RWV)
    if token._.is_rwp or token._.is_rwv:
        # 1.1.1 RWP rightmost item in the cluster
        if is_rightmost_in_cluster(token):

            # cope with ni => ne, vi => vă, and li => le
            interim_form = get_interim_form(token.lower_)

            # 1.1.1.1 obligatory sandhi to the right
            if r_nbor._.is_obligatory_host:
                surface_form = get_asyllabic_form(interim_form, "OBLIGATORY") + HYPHEN_TOKEN.text
                surface_forms.append(postverbal_hyphen+surface_form)

            # 1.1.1.2 no obligatory host to the right
            else:

                # 1.1.1.2.1 obligatory sandhi to the left for the u- and i-forms (lu ==> -l, mi ==> -mi)
                if l_nbor._.is_rwp:
                    # 1.1.1.2.1.1 e- or ă-forms
                    if ('e' in interim_form) or ('ă' in interim_form):
                        surface_form = interim_form
                        surface_forms.append(postverbal_hyphen+surface_form)
                        ## check for optional sandhi
                        if r_nbor._.vowel_initial_char and r_nbor._.vowel_initial and not \
                            r_nbor._.is_first_syllable_stressed and not r_nbor.lower_.startswith('o') and not\
                            r_nbor.lower_.startswith('e'):
                            surface_form = get_asyllabic_form(interim_form, "OPTIONAL") + HYPHEN_TOKEN.text
                            surface_forms.append(postverbal_hyphen+surface_form)
```

Figura 1: Example of RWP surface generation rules

| Feature | Preverbal forms | Postverbal forms | Example |
|---|---|---|---|
| RWP sequence length | L1 280 / L2 79 / L3 16 | L1 71 / L2 31 / L3 7 | Postverbal L2 / D - mi - le acum! / «Give them to me now!» |
| | Obligatory î-prothetic forms | Optional prothetic forms | |
| Î-prothetic contexts (Only preverbal L1) | 25 | î-prothetic 43 forms / no î-prothetic 50 forms | Optional î-prothetic form / De ce îmi dai mere? / «Why are you giving me apples?» |
| | Obligatory forms | Optional sandhi contexts | |
| Noî-prothetic contexts | 267 | RWP as syllabic host 32 / Context item as syllabic 50 | Obligatory RWP forms / De ce mi le dai? / «Why are you giving them to me?» |
| | Syllabic context item | Asyllabic context item | |
| Optional sandhi contexts with RWP as possible syllabic host | 15 | 17 | RWP as syllabic host / N-o v dbine. / «I don't see her/it well.» |
| | Syllabic RWP forms | Asyllabic RWP forms | |
| Optional sandhi contexts with context item as possible syllabic host | 25 | 25 | Context item as syllabic host / El ne-arat muzeul! / «He shows us the museum.» |

Tabela 2: Distribution of RWP features of target se...

put database is to establish a global granularity. These four functions [...] fair comparison and evaluation of... two phrasal syntactic a... based and GPT-based approaches... designed to simplify the v... textual features for generation RWP form.

## 6 Feature annotation and validation

The model relies on fine-grained information about both the current item and its context, which necessitate that the input data be linguistic... accordingly.

Finally, to ensure that hyph... a syllabicity, postverbality, ... ned according to official writing... the annotations are loaded... is_preverbal has been added to the... spaCy token features. In contrast... for processing the... position (ex 2) all items in... ding RWPs and auxiliary verbs... ted to the main verb by hyphens...

As already mentioned, ... entire set of 352 input sentences... use spaCy for several reasons:

Each Romanian input string is annotated... spaCy module, extended with func... ditional features; and the resulting... in JSON format. Task-relevant... are then manually corrected,... JSON structures since only task-... may contain errors. Since only task-... tations are manually correcte... tions may contain errors.

- it offers multi-language support for over 75 languages, including Romanian²;

- it offers a hybrid approach ... pports both machine learning models and linguistic rules;

- it is highly customizable and extensible (e.g., via CustomTokenizer ... as is_obligatory_host)

To address both obligatory and optional sandhi in RWP surface realization, the annotations have been extended with a set of functions that enrich token annotation with features for grapheme-phoneme disambiguation (vowel_final, vowel_ ...) and indicate whether an initial vowel is stressed (is_first_syllable_stressed).

Another set of functions has been devised to add high-level annotations, such as is_rwp, is_rightmost_in_cluster ...

## 7 Output testing

As with the development of any ... crafted grammar, the implementa... RWP surface generation module ... tive process involving careful ... finement, and extensive testing ... linguistic patterns were accu... Each iteration required revis... les, adding new ones, and addr... ses.

The spaCy-annotated, manual...

²https://spacy.io/models/ro

```
"ex021": {
  "ex021_input": "Arată ni muzeul!",
  "ex021_source": "Arată-ne muzeul!",
  "targets": {
    "ex021_t01": "Arată-ne muzeul!"
  }
}
```

Figura 2: Database entry with one target

```
"ex049": {
  "ex049_input": "De ce i dai mere?",
  "ex049_source": "De ce îi dai mere?",
  "targets": {
    "ex049_t01": "De ce îi dai mere?",
    "ex049_t02": "De ce-i dai mere?"
  }
}
```

Figura 3: Database entry

```
"ex242": {
  "ex242_input": "O împușcă în inimă.",
  "ex242_source": "O împușcă în inimă.",
  "targets": {
    "ex242_t01": "O împușcă în inimă.",
    "ex242_t02": "O-mpușcă în inimă.",
    "ex242_t03": "O împușcă-n inimă.",
    "ex242_t04": "O-mpușcă-n inimă."
  }
}
```

Figura 4: Database entry with four targets

```
"ex060": {
  "ex060_input": "Du te încolo!",
  "ex060_source": "Du-te încolo!",
  "targets": {
    "ex060_t01": "Du-te încolo!",
    "ex060_t02": "Du-te-ncolo!"
  }
}
```

Figura 5: Entry with optional targets

tput strings then comprises a... binations of these token surfa... of overgeneration are filtered... step of processing, as illustr... and 4. below.

```
<060> . . . . . . . . . . . . .
[4 tokens] Du te încolo!
[['Du'], ['-te', '-te-'],
 ['încolo', '-ncolo'], ['!']]
1. ('Du', '-te', 'încolo', '!')
2. ('Du', '-te', '-ncolo', '!')
3. ('Du', '-te-', 'încolo', '!')
4. ('Du', '-te-', '-ncolo', '!')
. . . . . . . . . . . . . . . . .
```

The remaining correct output... te încolo, Du-te-ncolo («Go away!»)... are checked against the target... ponding database entry in Figu...

mini-corpus based on the input data is loaded and used to re-analyze the input so that the linguistic features needed for... neration are available for the application of the rules. Finally, the output is checked against the corresponding set of targets (see Figures 2 – 4).

Since this implementation is primarily a proof of concept – for immediate context che- cking –, I kept the processing as simple as po- ssible. This simplicity means that the module has a stateless design, ... does not retain information... usly generated form of the current token and does not... of the previously generated... In standard Romanian or... can denote sandhi, postv... Given the lack of memory... hyphen may occasionally be...

For each token, a set of... nerated based on its context.

The hand-crafted grammar is... cifically for a defined input se... cular output set, making it ne... input linguistic choice of ex... tions more entirely error-free... or implementation. Neverthel... an attempt to faithfully imple... description outlined in Gers... with a computational-lingu...

## 8 Conclusions

In this article, I presented th... of a model for RWP surface rea... meaning that the module descri... Gerstenberger (2022). While... (2018) addresses only the rea... gator and... and do... ... implementation of the... study extends the approach by... both obligatory and optional... computational framework.

A further key contribution of... the compilation of a systemati...

linguistic annotations and implemented constraints. The implementation serves not only as a proof of concept but also as a robust instrument for validating the model, demonstrating its practical applicability, and ensuring that its predictions align with empirical data. Due to its transparency and proximity to surface forms, the model is adaptable to any constraint-based linguistic framework.

Since the resources used in this study, specifically, the input test data from the surface generation module – were created by us. However, they may contain errors. These resources are freely available at https://github.com/ciprian-NO/rbl... for the research community to use and expand.

# References

Roxana-Maria Barbu and Ida Toivonen. Romanian Object Clitics: Grammatical agreement and ... Proceedings of the LFG18 Conference, pages 67-87. Stanford. CSLI Publications.

Kenneth R. Beesley and Lauri Karttunen. Finite-State Morphology. CSLI Publications, Stanford, CA.

Emily M. Bender and D. Terence Langendoen. 2010. Computational linguistics in support of linguistic theory. Linguistic Issues in Language Technology, 3.

Eulàlia Bonet. 1994. The person-case constraint: A morphological approach. MIT Working Papers in Linguistics, 22:33-52.

Anca Chereches. 2014. A Prosodic ... Romanian Pronominal Clitics. University of Pennsylvania Working Papers in Linguistics, volume 20. Available at https://repository.upenn.edu/pwpl/vol20...

Noam Chomsky and Morris Halle. 1968. The Sound Pattern of English. Harper & Row.

Ann A. Copestake. 2001. Implementing typed feature structure grammars. CSLI lecture notes series.

Carmen Dobrovie-Sorin. 1999. ... gories: The case of Romanian. ... Riemsdijk, editor, in the Languages of Europe. Mouton de Gruyter.

Carmen Dobrovie-Sorin and Ion Giurgea, editors. 2013. A Reference Grammar of Romanian. Volume 1: The noun phrase. John Benjamins.

Berthold Crysmann and Tracy Holloway King. Clitics in ... Phonology, Morphology ... John Benjamins.

... Gerstenberger. 2007. ... realization ... Proceedings of EUROLAN 2007, University of Iasi, Romania. at https://www.researchgate.net/... A morphology-based realization model for generation.

Ciprian-Virgil Gerstenberger. ... Romanian weak pronoun generation. Languages at the Crossroads: ... creditation and Context of Use ... of the 35th Edition of the International Conference of the Spanish Association ... pages 215-226. University ... Available at https://www.researchgate.net/publication/332671041_A_Grammar...manian_weak_pronoun_generati...

Ciprian-Virgil Gerstenberger. are romanian clitics pronouns? Compound nouns? 47(1):37-57. 2018.

Bruce Hayes, Bruce Tesar, and Kie Zuraw. OTSoft 2.5 [software package]. linguistics.ucla.edu/people/hayes/...

Ronald M. Kaplan, John T. Maxwell, Tracy Holloway King, and Richard Crouch. ... Integrating finite-state technology with ... grammars. Proceedings of the ESSLLI Workshop on Combining Shallow and Deep Processing for NLP.

Udo-Michael Klein. 2007. ... of argument structure in Romanian and ... PhD thesis, University of London. at http://hdl.handle.net/11858/00-... 0012-8EEC-D.

Kimmo Koskenniemi. 1983. Two-level morphology ... A, Volume 83, pages 683+685.

Géraldine Legendre. 2001. Positioning verbal clitics at PF. In (Gerlach and ... 2001).

Paola Monachesi. 2001. Clitic p... Romanian verb complex. In (Gerlach and van Riemsdijk ... 2001).

Paola Monachesi. 2005. Complex in Romanian: A case study in Grammar ... Oxford University Press.

John J. Ohala. 2008. The emergent syllable in speech. Taylor & Francis ...

OpenAI. 2024. ChatGPT [Large La...] https://chat.openai.com

Alexandra Popescu. 2000. The morphophonology of the Romanian clitic sequence. In Lingua, volume 110, pages 773–799. Elsevier.

Alexandra Popescu. 2003. Morphophonologische Phänomene des Rumänischen. PhD thesis, University of Düsseldorf. Available at https://docserv.uni-duesseldorf.de/servlets/-DocumentServlet?id=3187.

Kan Sasaki and Daniela C Iuianu. 2000. An Optimality Theoretic Account for the Distribution of Pronominal Clitics in Romanian. Technical report, University of Tsukuba. Available at https://www.academia.edu/33668351/An_optimality_theoretic_account_for_the_distribution_of_pronominal_clitics_in_Romanian.

Stanca Somesfalean. 2007. Form and Interpretation of Clitics. PhD thesis, Université du Québec à Montréal. Available at https://archipel.uqam.ca/9628/.

Oana Savescu Ciucivara. 2010. A Syntactic Analysis of Pronominal Clitic Clusters in Romance - The view from Romanian. PhD thesis, New York University. Available at https://www.proquest.com/docview/304954033.