

# Contemporary Persian Inflectional Analyzer

## Davood Heidarpour

MSc in Computational Linguistics; Faculty of New Sciences and Technologies; University of Tehran; Tehran, Iran;  
Email: d.heidarpour@ut.ac.ir

## Elham S. Sebt

MSc in Computational Linguistics; Faculty of New Sciences and Technologies; University of Tehran; Tehran, Iran;  
Email: e.sebt@ut.ac.ir

## Mahmood Bijankhan

PhD in General Linguistics; Full Professor; Faculty of Literature and Humanities; University of Tehran; Tehran, Iran;  
Email: mbjkh@ut.ac.ir

## Mostafa Salehi\*

PhD in Computer Engineering; Associate Professor;  
Faculty of New Sciences and Technologies; University of Tehran; Tehran, Iran Email: mostafa\_salehi@ut.ac.ir

## Hadi Veisi

PhD in Computer Engineering; Assistant Professor;  
Faculty of New Sciences and Technologies; University of Tehran; Tehran, Iran Email: h.veisi@ut.ac.ir

Iranian Journal of  
**Information  
Processing and  
Management**

Iranian Research Institute

for Information Science and Technology  
(IranDoc)

ISSN 2251-8223

eISSN 2251-8231

Indexed by SCOPUS, ISC, & LISTA

Vol. 36 | No. 4 | pp. 945-970

Summer 2021

<https://doi.org/10.52547/ijpm.36.4.945>



Received: 19, Sep. 2019 | Accepted: 07, Mar. 2021

**Abstract:** In recent years, the use of informal writing in Persian has grown significantly due to the increasing expansion of cyberspace and social media and platforms, and the tendency of users to bring the written language closer to colloquial speech. But on the other hand, proper tools to process this language register are not developed very much. One of the tools for low level processing of textual data is an inflectional analyzer. However, such tools are not developed for this register yet. Informal words have their own structures, stems, morphemes and clitics and they also make use of formal structures and units. Moreover, this register also consists of formal words so any analyzer for informal words should have the potential to analyze formal words, too. In this paper, it is tried to cover all inflectional structures of informal Persian language to build an inflectional analyzer. A corpus of most of its known sub-registers is constructed to extract words, morphemes and inflectional rules and morphotactics. A part of this corpus is used for testing the analyzer. After extracting 1786 unique words of the test part, inflectional analyzer f-measure is equal to 97.67%.

\* Corresponding Author

This tool can be used in computational processing of Persian language and it can also be used in teaching Persian, specifically colloquial Persian to non-Persian learners.

**Keywords:** Computational Linguistics, Inflectional Analyzer, Informal Persian Register, Contemporary Persian, FST, Persian Teaching

# تحلیلگر تصریفی فارسی معاصر

داود حیدرپور

کارشناسی ارشد زبان‌شناسی رایانشی؛ دانشکده علوم و فنون نوین؛ دانشگاه تهران؛ تهران، ایران؛  
d.heidarpour@ut.ac.ir

الهام‌سادات سبٹ

کارشناسی ارشد زبان‌شناسی رایانشی؛ دانشکده علوم و فنون نوین؛ دانشگاه تهران؛ تهران، ایران؛  
e.sebt@ut.ac.ir

محمود بی‌جن خان

دکتری زبان‌شناسی؛ استاد تمام؛ گروه زبان‌شناسی؛ دانشگاه تهران؛ تهران، ایران  
mbjkan@ut.ac.ir

مصطفی صالحی

دکتری مهندسی کامپیوتر؛ دانشیار؛ دانشکده علوم و فنون نوین؛ دانشگاه تهران؛ تهران، ایران؛  
mostafa\_salehi@ut.ac.ir

هادی ویسی

دکتری مهندسی کامپیوتر؛ استادیار؛ دانشکده علوم و فنون نوین؛ دانشگاه تهران؛ تهران، ایران؛  
h.veisi@ut.ac.ir



دریافت: ۱۳۹۸/۰۶/۲۸ | پذیرش: ۱۳۹۹/۱۲/۱۷ | مقاله برای اصلاح به مدت ۴ ماه نزد پدیدآوران بوده است.

نشریه علمی | رتبه بین‌المللی  
پژوهشگاه علوم و فناوری اطلاعات ایران  
(ایرانداک)

شاپا (چاپی) ۲۲۵۱-۸۲۲۳

شاپا (الکترونیکی) ۸۲۳۱-۲۲۵۱

نمایه در SCOPUS، ISI، LISTA،

jipm.irandoc.ac.ir

دوره ۳۶ | شماره ۴ | صص ۹۴۵-۹۷۰

تابستان ۱۴۰۰

<https://doi.org/10.52547/jipm.36.4.945>



چکیده: در سال‌های اخیر، کاربرد گونه نوشتاری غیررسمی زبان فارسی به دلیل گسترش روزافزون فضای مجازی و شبکه‌های اجتماعی و تمایل کاربران به نزدیک کردن زبان نوشتار به گفتار رشد چشمگیری داشته است. با وجود این، ابزارهای پردازش این گونه زبانی به میزان لازم توسعه داده نشده است. تحلیلگرهای تصریفی از جمله ابزارهایی است که در پردازش زبانی کاربرد وسیعی دارد و تاکنون برای گونه غیررسمی طراحی و پیاده‌سازی نشده است. با توجه به این که گونه نوشتاری غیررسمی در کنار واژگان و قواعد صرفی و نحوی مختص به خود، در بخشی از واژه‌ها و ساختارها با گونه رسمی مشترک است، در این پژوهش با پوشش فارسی رسمی و غیررسمی اولین ابزار تحلیل تصریفی فارسی معاصر برای همه اقسام واژه توسعه داده شده و تلاش شده همه ساختارهای تصریفی واژه‌های فارسی غیررسمی پوشش داده شود. این ابزار به صورت قاعده‌مند و مستقل از بافت و با بهره‌گیری از مبدل حالت محدود، پی‌بست‌ها و وندهای تصریفی رسمی و غیررسمی را در واژه‌های زبان شناسایی و تحلیل کرده، ستاک‌های رسمی و غیررسمی را نیز استخراج می‌کند. به‌منظور پوشش دادن تمام ساختارها و حالت‌های تصریفی، با توجه به

رویکرد مستقل از بافت، الگوریتم برای هر واژه، تمام خوانش‌ها و معانی گوناگونی را که می‌تواند بسته به قرارگیری در بافت‌های گوناگون داشته باشد، تحلیل و ارائه می‌کند. به‌منظور استخراج و بررسی واژگان و قواعد تصریفی و نگارشی گونه غیر رسمی، پیکره فارسی معاصر از سیاق‌ها و زیرسیاق‌های گوناگون این گونه زبانی تهیه شده و در طراحی و آزمون تحلیلگر مورد استفاده قرار گرفت. آزمون تحلیلگر با استفاده از ۱۷۸۶ واژه یکنای استخراج شده از پیکره، نتیجه ۹۶/۶۷ درصد را در معیار اف به‌دست داده است. از این ابزار می‌توان در انواع تحلیل‌ها و کاربردهای پردازش رایانه‌ای زبان فارسی و همچنین در آموزش فارسی، به‌ویژه محاوره فارسی به غیر فارسی‌زبانان استفاده کرد.

**کلیدواژه‌ها:** پردازش زبانی، تحلیلگر تصریفی، گونه غیررسمی فارسی، فارسی معاصر، مبدل حالت محدود، آموزش فارسی

## ۱. مقدمه

گونه‌ای از زبان فارسی که به‌عنوان زبان رسمی کشور در اسناد و مکتوبات حکومتی، آموزشی، علمی و نیز روزنامه‌ها و نشریات به کار می‌رود، گونه نوشتاری معیار یا رسمی است. در کنار این گونه می‌توان به گونه غیررسمی اشاره کرد که به‌عنوان نوشتار ارتباطی در شبکه‌های اجتماعی، پیام‌رسان‌های تلفنی، وبلاگ‌ها و گاهی سایت‌ها و نامه‌های الکترونیکی استفاده می‌شود و کاربر تلاش می‌کند زبان نوشتار را هرچه بیشتر به گونه گفتاری نزدیک کند. هرچند این گونه نوشتاری، ساختاری مختص به خود دارد و از لحاظ صرفی، نحوی و واژگانی با زبان نوشتاری معیار متفاوت است، اما در کنار تفاوت‌ها وجوه اشتراک دو گونه نیز بالاست. در بررسی آماری پیکره‌ای که برای این پژوهش جمع‌آوری شده (شامل بیش از ۴۹ هزار واژه که به‌صورت دستی تقطیع<sup>۱</sup> شده)، ۸۳/۵ درصد (بیش از ۴۱ هزار) واژه‌های متون فارسی غیررسمی با زبان فارسی گونه رسمی مشترک است و تنها ۱۶/۵ درصد واژه‌ها غیررسمی است (بخش ۳). لازم به ذکر است که واژه غیررسمی بر اساس این پژوهش، واژه‌ای است که حداقل یک عنصر غیررسمی در ساختار خود داشته باشد و این عنصر ریشه یا هر یک از اجزای تصریفی واژه است.

در حوزه پردازش رایانه‌ای زبان، تحلیل صرفی به‌عنوان پردازش پایه در بسیاری از پردازش‌های زبانی سطح بالاتر کاربرد دارد، اما ابزارهای شناسایی فارسی غیررسمی بر خلاف گسترش بسیار این گونه زبانی در رسانه‌های اجتماعی و فضای وب، توسعه بسیار

1. tokenize

اندکی یافته‌اند.

پژوهش‌های انجام‌شده در حوزه ساخت‌واژه زبان فارسی بیشتر به گونه رسمی زبان پرداخته‌اند، از جمله، «اسلامی» و همکاران (۱۳۸۳)، «اسلامی و علی‌زاده لجمیری» (۱۳۸۸)، Shamsfard, Jafari and Ilbeygi (2010) و «مواجی، اسلامی و وزیرنژاد» (۱۳۹۰). تحلیلگرهای تصریفی که در این پژوهش‌ها پیاده‌سازی شده، مختص زبان فارسی رسمی است و قادر به تصریف فارسی غیررسمی نیست. این پژوهش‌ها بر اساس ساختار تصریفی ارائه شده توسط «اسلامی» و همکاران (۱۳۸۳) طراحی شده‌اند. تنها پژوهش Megerdoomian (2000) بر اساس ساختار تصریفی‌ای که خود او استخراج و تبیین کرده، پیاده‌سازی شده است. تعداد اندکی از پژوهش‌ها نیز به فارسی غیررسمی پرداخته‌اند. «مگردومیان» به صورت پراکنده ساختارهای غیررسمی را شناسایی کرده است (Megerdoomian 2008). این ساختارها می‌بایست به تحلیلگر رسمی او -ساختارهایی که خود او استخراج و تبیین کرده- افزوده می‌شد تا تحلیلگر تصریفی بتواند در کنار تصریف واژگان رسمی، واژگان غیررسمی را نیز پوشش دهد، اما این طرح ناتمام باقی مانده است. گرچه این پژوهش قبل از فراگیر شدن پیام‌رسان‌ها و شبکه‌های اجتماعی انجام شده، اما با توجه به رشد وبلاگ‌های فارسی و گسترش کاربرد فارسی غیررسمی در آن اولین پژوهشی است که به فارسی غیررسمی در فضای مجازی پرداخته است. «تازه‌جانی و بحرانی» با انجام پژوهشی محدود، تنها بخش کوچکی از افعال غیررسمی را پوشش داده و تحلیلگری تصریفی برای آن ساخته‌اند (۱۳۹۲). «طیب‌زاده» نیز مانند «مگردومیان» با مطالعه‌ای پیکره‌بنیاد، اما برخلاف آن محدود به داستان‌ها و نمایشنامه‌های صد سال اخیر فارسی، رخدادهای کلمات غیررسمی را بررسی کرده و ساختار مشاهده شده در کلمات را گزارش کرده است (۱۳۹۸). این پژوهش صرفاً نظری است و فاقد ابزار تحلیل رایانه‌ای این گونه کلمات است.

استفاده از زبان غیررسمی در اینترنت در میان سخنوران سایر زبان‌ها نیز رایج بوده و تلاش برای تحلیل رایانه‌ای آن نیز در جریان است. برای مثال، Torjmen and Haddar (2018) برای تحلیل گونه بیانی تونسسی و Alshargi and Rambow (2016) برای تحلیل گونه بیانی یمنی (صناعی) و مراکشی، که همگی متفاوت از عربی استاندارد مدرن (Ryding 2005) یا عربی نوشتاری مدرن<sup>۱</sup> (Badawi, Carter and Gully 2013) هستند، با شناسایی

1. Modern Written Arabic (MWA)

ساختار چنین کلماتی (با استخراج از پیکره‌هایی که برای این منظور ساخته‌اند) و افزودن آن‌ها به ساختار مبدل‌های<sup>۱</sup> عربی استاندارد یا نوشتاری مدرن، تحلیلگرهای تصریفی برای گونه‌های بیانی خود طراحی کرده‌اند.

با توجه به کاستی‌های این حوزه، در پژوهش حاضر تلاش شده است که ساختارهای تصریفی فارسی غیررسمی به‌طور کامل بررسی و شناسایی شده، و تحلیلگر تصریفی فارسی غیررسمی بر مبنای قواعد جامع و مانع در این حوزه طراحی و پیاده‌سازی شود. هدف در طراحی این تحلیلگر شناسایی و تحلیل ستاک، وندهای تصریفی و پی‌بست‌های واژه است که در سامانه‌های پردازش زبان طبیعی<sup>۲</sup> کاربرد دارد.

از آنجا که در گونه غیررسمی از گونه رسمی نیز استفاده می‌شود، بنابراین، هر تحلیلگر واژه برای فارسی غیررسمی می‌بایست قادر به تحلیل واژه‌های فارسی رسمی هم باشد. فارسی‌ای که در این پژوهش مدنظر است، فارسی امروزی است که در نوشتار و گفتار فارسی‌زبانان در گونه رسمی و غیررسمی استفاده می‌شود و شامل گونه ادبی و تاریخی نمی‌شود. بنا بر نظر «لازار» فارسی معاصر کاربرد امروزی فارسی را در هر شکلش در نظر دارد و نه ساختاری که لزوماً در متون کلاسیک فارسی وجود دارد و یا تعریف شده است (Lazard 1992).

## ۲. چارچوب پیشنهادی

رویکرد رایج در تحلیل زبانی قاعده‌مند و در سطح واژه، استفاده از مبدل حالت محدود<sup>۳</sup> است که می‌تواند دو رشته را به یکدیگر نگاشت کند. این نگاشت بین واژه و تمام قطعات سازنده قاعده کلی واژه صورت می‌گیرد. بنابراین، ورودی مبدل می‌تواند واژه و خروجی آن قاعده سازنده واژه باشد. به‌منظور پیاده‌سازی یک تحلیلگر با استفاده از مبدل حالت محدود به سه منبع زبانی نیاز داریم: واژگان، قواعد تصریفی و قواعد نگارشی (Jurafsky and Martin 2008).

واژگان<sup>۴</sup> فهرستی از ستاک‌ها و وندهای زبان همراه با اطلاعاتی مثل مقوله دستوری واژه است. قواعد هم‌نشینی<sup>۵</sup> تک‌واژه‌ها، مدلی از هم‌نشینی و باهم‌آیی تک‌واژه در یک واژه

1. transducer

۲. نظیر تحلیل معنایی، نظر کاوی، بازشناسی گفتار و غیره.

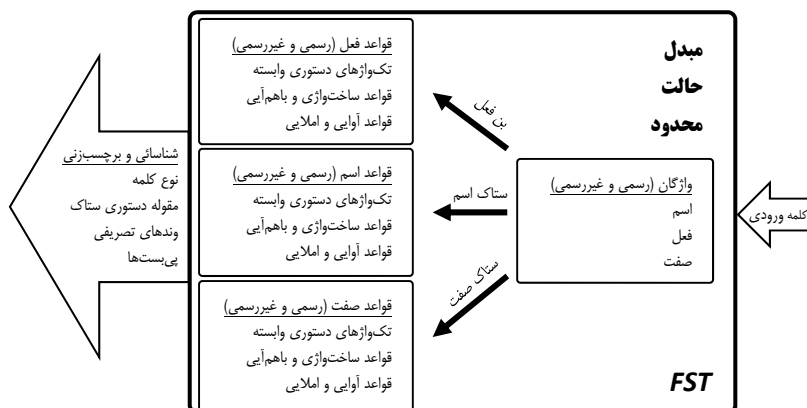
3. finite state transducer (FST)

4. lexicon

5. morphotactic rules

ارائه می‌کند و قواعد املایی، تغییرات املایی ناشی از تغییرات واژ-واجی<sup>۱</sup> را هنگام ترکیب دو تک‌واژ توصیف می‌کند.

در این پژوهش به منظور پردازش تصریفی قاعده‌مند از ابزار «فوما»<sup>۲</sup> (Hulden 2009)، که یک ابزار مبدل حالت محدود متن آزاد و رایگان است، استفاده شده است. «فوما» این امکان را فراهم می‌آورد که با قرار دادن چندین مبدل در سیستم، واژه ورودی را از یکی یا از همه آن‌ها عبور داد و قاعده تولید کرد. هر مبدل می‌تواند برای تغییر و یا شناسایی بخشی از قواعد ساخت واژه به کار رود. همچنین، کاربر می‌تواند با استفاده از کتابخانه‌های تعریف شده یا به صورت شخصی شده در مبدل‌ها جست‌وجو انجام دهد و یا ترتیب و نحوه اجرا شدن مبدل‌ها را شخصی‌سازی کند. شکل ۱، ساختار مفهومی یک مبدل را نشان می‌دهد. کلمه ورودی پس از شناسایی ستاک و ساختار تصریفی و برجسب‌زنی همه اجزای سازنده آن در خروجی نمایش داده می‌شود.



شکل ۱. مدل مفهومی مبدل حالت محدود تحلیلگر تصریفی

## ۲-۱. چالش‌های پیاده‌سازی تحلیلگر تصریفی فارسی معاصر

چنان‌که ذکر شد، گونه‌ای از زبان که به‌عنوان فارسی معاصر در این پژوهش مورد توجه است، گونه‌ای ارتباطی است که به شکل گسترده در شبکه‌های اجتماعی، پیام‌رسان‌ها و وبلاگ‌ها مورد استفاده قرار می‌گیرد و دارای ویژگی‌های ارتباطی نزدیک

به گفتار است. با توجه به محدودیت‌های فنی رسانه<sup>۱</sup> نوشتاری فارسی و فقدان استاندارد، کاربران از شیوه‌های گوناگونی برای نگارش استفاده می‌کنند. اما به‌رغم تغییرات آوایی و یکدست نبودن قوانین نگارشی، فارسی غیررسمی در بسیاری از موارد در ساختار تصنیفی خود از قواعدی تبعیت می‌کند که در محدود پژوهش‌های این حوزه به برخی از آن‌ها اشاره شده است. اما لازم است برای تحلیل دقیق و درست، با بررسی داده‌های استاندارد این گونه‌ی زبانی، ساختار جامع و مانعی برای آن تعریف شود. با توجه به در دسترس نبودن داده‌ی استاندارد برای این گونه‌ی نوشتاری، نخستین قدم فراهم آوردن داده‌ی مناسب به‌منظور بررسی و استخراج الگوهای نوشتاری، تک‌واژه‌های تغییر یافته و قواعد ساخت‌واژی است. با توجه به این که در این گونه‌ی نوشتاری به‌طور معمول واژه‌ها چه در شکل رسمی و چه در شکل غیررسمی، تابع الگوی نگارشی یکسان و ثابتی نیستند و پوشش و ذخیره‌ی همه‌ی این واژه‌ها در پیاده‌سازی تحلیلگر امکان‌پذیر نیست، بنابراین، بخشی از واژه‌ها، به‌ویژه آن‌هایی که تحت تأثیر تغییرات آوایی قرار می‌گیرند، تبدیل به واژه‌های خارج از واژگان تحلیلگر می‌شوند و شناسایی آن‌ها به‌راحتی امکان‌پذیر نیست. برای به حداقل رساندن تأثیر چنین مواردی در دقت نهایی تحلیلگر لازم است به نکات زیر توجه شود:

◇ حذف برخی حروف: همزه در پایان واژه‌ها یا بر روی پایه‌های «و»، «ا»، و «ی» وجودی غیر قطعی دارد و ممکن است حذف شود؛ مثلاً «أعضاء: اعضا». حرف «ه» در نقش واکه در صورتی که در پایان یک تک‌واژه قرار بگیرد، در بسیاری از موارد، در هنگام اتصال به تک‌واژه‌ی دیگر حذف می‌شود؛ مثلاً «به‌طوریکه: بطوریکه، همه شون: همشون».

◇ فاصله‌گذاری: فاصله، نیم‌فاصله و اتصال، سه حالت در کنار هم قرار گرفتن واژه‌ها و یا اجزای واژه‌هاست. گرچه «فرهنگستان زبان و ادب فارسی» استاندارد برای آن در نظر گرفته (صادقی و زندی مقدم ۱۳۸۵)، اما به‌طور معمول، در نگارش فارسی رعایت نمی‌شود؛ مثل «کتاب‌ها، کتابها و کتاب‌ها» و در نگارش غیررسمی نیز تنوع بیشتری تولید می‌کند.

◇ خطای املائی: خطای املائی یکی از مشکلات اجتناب‌ناپذیر در همه‌ی پردازش‌های متنی به‌ویژه فارسی غیررسمی است. بخشی از این خطاها که منشأ نگارشی دارد، با

1. medium



برطرف کردن سایر چالش‌های پردازش فارسی غیررسمی قابل کنترل است؛ مانند جایگزینی حروف هم‌صدا با یکدیگر. سایر خطاهای ناشی از حروف چینی نیز با روش‌های شناسایی و رفع خطای املائی می‌بایست برطرف شود.

### ۳. دادگان

در تهیه دادگان برای این پژوهش تلاش شد تمام زیرسیاق‌های فارسی غیررسمی شناسایی شده، برای هر کدام متونی تهیه، بررسی و برجسب‌گذاری شود. مجموعه پیکره جمع‌آوری شده نزدیک به پنجاه‌هزار واژه است. واژگان این پیکره به شکل دستی تقطیع شده و حدود ۳۰ درصد آن برای ارزیابی کنار گذاشته شد. از بقیه برای بررسی ساخت‌های تصریفی، آمارگیری از میزان توزیع ساخت‌های مختلف، بررسی و استخراج واژگان و ساخت‌های غیررسمی استفاده شد. تمام واژه‌های غیررسمی این پیکره به شکل دستی معادل‌نویسی شده و بن‌واژه<sup>۱</sup> آن‌ها نیز استخراج شده است. ساختار پیکره‌ای که به‌طور متوازن از این زیرسیاق‌ها (که همگی مربوط به زبان فارسی غیررسمی است) استفاده می‌کند، از نوع پیکره نمونه‌گیری شده<sup>۲</sup> است. برای هر سیاق<sup>۳</sup> زبان فارسی غیررسمی نمونه‌ها، هم در مرحله جمع‌آوری اولیه و هم در مرحله گزینش به شکل تصادفی<sup>۴</sup> انتخاب شده است (McEnry and Hardie 2011). حجم داده هر بخش هم به نسبت کاربرد آن زیرسیاق انتخاب شده است. جدول ۱، داده‌های جمع‌آوری شده در تهیه پیکره فارسی معاصر را به تفکیک سیاق‌ها و زیرسیاق‌ها نشان می‌دهد.

### ۴. قواعد و واژگان فعلی

در این بخش واژگان فعلی تحلیلگر که شامل بن‌افعال و تک‌واژه‌های وابسته متصل به فعل است و همچنین، ساخت تصریفی افعال در گونه غیررسمی توضیح داده می‌شود.

1. lemma

2. sampling corpus

3. register

4. random sampling

### جدول ۱. سیاق‌ها و زیرسیاق‌های موجود در پیکره فارسی معاصر

ردیف	سیاق	زیرسیاق و حجم کلی
۱	گفت‌وگو و بیان واقعی <sup>۱</sup>	مصاحبه، سخنرانی از پیش آماده‌شده، اجرای رادیو/ تلویزیونی از روی متن؛ ۱۲۴۸ واژه
۲	گفت‌وگو و بیان خیالی <sup>۲</sup>	رمان، نمایش‌نامه، فیلم‌نامه، زیرنویس فیلم، شعر محاوره؛ ۷۵۶۷ واژه
۳	نظر کاربران فضای مجازی	مصرف‌کننده کالا یا خدمات، مقالات خبری و متفرقه؛ ۹۰۵۰ واژه
۴	پیام‌رسان تلفنی	مکالمه شخصی، مکالمه گروهی، اطلاع‌رسانی؛ ۱۱۹۴۲ واژه
۵	شبکه اجتماعی	پست‌ها، پاسخ‌ها؛ ۸۴۳۱ واژه
۶	وبلاگ	پست‌ها، پاسخ‌ها؛ ۸۹۰۳ واژه
۷	نامه شخصی و خاطرات روزانه؛ ۲۱۴۶ واژه	

#### ۴-۱. واژگان فعلی

در ساختمان فعل‌های رسمی ساده و پیشوندی از بن‌های فعل ماضی و مضارع استفاده می‌شود. بن‌های رسمی از واژگان «زایا» (اسلامی و همکاران ۱۳۸۳)، «ویراستیار» (کاشفی، نصری و کنعانی ۱۳۸۹)، واژگان «پرلکس» (Sagot and Walther 2010)، و پیکره جمع‌آوری‌شده فارسی غیررسمی به‌دست آمده است. بن فعل‌های حال و گذشته غیررسمی نیز از پیکره جمع‌آوری‌شده و هم با معادل‌سازی بن فعل‌های رسمی به‌دست آمده است که در مجموع، شامل ۵۱۶ بن ماضی ساده رسمی و ۸۷ بن غیررسمی، ۴۹ بن ماضی پیشوندی رسمی و ۲۳ بن غیررسمی، ۳۷۹ بن مضارع ساده رسمی و ۸۷ بن غیررسمی و نیز ۴۹ بن مضارع پیشوندی رسمی و ۲۲ بن غیررسمی است. شناسه‌ها و پی‌بست‌ها با توجه به پژوهش‌های انجام‌شده، مانند (Megerdooian 2008)، «اسلامی» و همکاران (۱۳۸۳)، «انوری و گیوی» (۱۳۹۱)، «بی‌جن‌خان» (۱۳۸۶) و از پیکره جمع‌آوری‌شده به‌دست آمده است (جدول ۲).

در ساخت افعال گونه غیررسمی استفاده از پی‌بست‌های فاعلی و مفعولی رایج است. پی‌بست‌های مفعولی «م، ت، ش، مان، تان، شان» در متون رسمی نیز به‌ندرت استفاده می‌شود، اما در متون غیررسمی صیغه‌های جمع بیشتر به‌صورت «مون، تون، شون» استفاده می‌شود. پی‌بست فاعلی «ش» مختص گونه محاوره است و در فعل‌های سوم‌شخص به‌کار

1. non-fictional

2. fictional

می‌رود. بیشتر نموده‌های فعل ماضی مانند ماضی ساده، استمراری و بعید، صرف نظر از گذرا یا ناگذر بودن فعل پی‌بست فاعلی را می‌پذیرند. پی‌بست فاعلی می‌تواند به فعل‌های مضارع اخباری و التزامی لازم نیز افزوده شود.

### جدول ۲. شناسه‌های فارسی رسمی و غیررسمی<sup>۱</sup>

شخص و شمار	بعد از همخوان‌ها و واکه‌ی		بعد از واکه‌آ		بعد از واکه و	
	رسمی	غ.رسمی	رسمی	غ.رسمی	رسمی	غ.رسمی
۱ مفرد	م	می‌خندم، می‌زیم، می‌گم	م (ی)م	م	م (ی)م	م (ی)م
۲ مفرد	ی	می‌خندی، می‌زیی، می‌گی	ی (ی)ی	ی	ی (ی)ی	ی (ی)ی
۳ مفرد	مضارع: اد ماضی: Ø	مضارع: اه ماضی: Ø می‌گه، می‌داشت	د (ی)د	د	د (ی)د	د (ی)د
۱ جمع	یم	می‌خندیم، می‌زییم، می‌گیم	یم (ی)یم	یم	یم (ی)یم	یم (ی)یم
۲ جمع	ید	می‌خندید، می‌زیید، می‌گین	ید (ی)ید	ید / این / این (ی)ین	ید (ی)ید	ید (ی)ید
۳ جمع	ند	می‌خندند	ند (ی)ند	ند	ند (ی)ند	ند (ی)ند

- (۱). اومدش و بچه‌ها رو به صف کرد، بعدشم بردشون بیرون.
  - (۲). بچه‌های ما همش سراغ دختر کوچولو تونو می‌گیرن، فردا که میایین با خودتون بیارینش.
  - (۳). خودش گفت می‌خواسته ماجرا رو بگه ولی اونا نداشتن.
- افزون بر پی‌بست‌های فاعلی و مفعولی پی‌بست‌های «و»، «م» نیز که صورت کوتاه‌شدهٔ حرف عطف «و» و حرف ربط «هم» است و نیز پی‌بست تأکید «ها» می‌توانند به ساخت افعال گونهٔ غیررسمی اضافه شوند (Megerdooian 2008؛ شقاقی ۱۳۹۴) که در توزیع تکمیلی یکدیگر هستند.
- (۴). فایده نداشت اگه می‌گفتیم گوش نمی‌کرد.

۱.!: بدون هیچ فاصله‌ای به تک‌واژ قبل متصل می‌شود. ×: با نیم فاصله به تک‌واژ قبل متصل می‌شود.

(۵). اونم رفتو با یه مامور برگشت.

همچنین، در نوشتار گونه غیررسمی از پی‌بست تأکیدی «ا!» نیز استفاده می‌شود که با بررسی پیکره جمع‌آوری شده شناسایی و استخراج شد.

(۶). اگه حرف گوش نکنی نمی‌برمتا!

#### ۴-۲. ساخت تصریف فعلی

به‌منظور طراحی و پیاده‌سازی تحلیلگر تصریفی فارسی معاصر قواعد مربوط به ساخت تصریفی افعال گونه رسمی از منابع و پژوهش‌های انجام‌شده در این حوزه از جمله از «اسلامی» و همکاران (۱۳۸۳) و «انوری و گیوی» (۱۳۹۱) استخراج شد و مورد استفاده قرار گرفت. ساخت تصریفی فعل‌ها در گونه غیررسمی نیز از الگوی ساخت فعل در گونه رسمی تبعیت می‌کند، با این تفاوت که افعال غیررسمی امکان پذیرش طیف گسترده‌تری از پی‌بست‌ها را دارند که به‌عنوان قواعد افزوده به تحلیلگر ساخت‌واژی معرفی شده است. در ساخت فعل غیررسمی حداقل یک عنصر غیررسمی باید وجود داشته باشد که عبارت است از: بن غیررسمی فعل، شناسه غیررسمی، پی‌بست مفعولی، پی‌بست فاعلی، تک‌واژه‌های محاوره (جدول ۶). به‌عنوان مثال، فعل «گفتیشان» گرچه از بن رسمی «گفت»، شناسه رسمی «ی» و پی‌بست ضمیری رسمی «شان» تشکیل شده است، اما دارای ساختی غیررسمی است؛ زیرا فعل‌های رسمی پی‌بست مفعولی نمی‌پذیرند (Megerdooonian 2000). فعل «رفتش» نیز که از بن رسمی و شناسه سوم شخص مفرد (Ø) تشکیل شده، به این دلیل که پی‌بست فاعلی «ش» پذیرفته است، یک فعل غیررسمی به حساب می‌آید (اسلامی و علی‌زاده لجمیری ۱۳۸۸). در این پژوهش با بررسی تحقیقی که به‌صورت ناتمام در این زمینه توسط Megerdooonian (2008) انجام شده، همین‌طور، منابع مشترک با فعل رسمی و دیگری پیکره جمع‌آوری شده، الگوی تصریفی افعال غیررسمی مطابق شکل ۲، استخراج و مورد استفاده قرار گرفت.

$$\left( \text{اجزای ساختمان غیررسمی} \right) + \left[ \begin{array}{l} \left( \text{پی‌بست مفعولی} \right) \\ \left( \text{پی‌بست فاعلی} \right) \end{array} \right] + \text{ساختمان فعل رسمی}$$

شکل ۲. تصریف افعال غیررسمی

فعل‌های گذرا می‌توانند پی‌بست مفعولی و یا پی‌بست فاعلی بپذیرند. در مقابل، فعل‌های ناگذر تنها می‌توانند پی‌بست فاعلی بپذیرند. در فعل‌های گذرا پذیرفتن پی‌بست فاعلی یا پی‌بست مفعولی سوم شخص مفرد هر دو منجر به افزوده شدن «ش» به انتهای بن فعل می‌شود که ابهام ایجاد می‌کند.

#### ۴-۲-۱. موارد استثنا در ساخت فعل امر

در چارچوب قاعده تصرفی فعل امر غیررسمی می‌توان چند مورد استثنا شناسایی و معرفی کرد:

◇ در ساخت فعل‌های غیررسمی این امکان وجود دارد که بن‌های رسمی یا غیررسمی با شناسه‌های رسمی یا غیررسمی ترکیب شوند؛ به‌عنوان مثال: «بگذرانید، بگذرانین، بگذروید، بگذرونین». در ساخت فعل امر مفرد در صورتی که بن مضارع مختوم به واکه (رو، گو، شو، ده) باشد، از بن غیررسمی استفاده نمی‌شود.

(۷) مصدر گفتن - بن رسمی: گو «بگوئید»، بن غیررسمی: گگ «بگید» ساخت امر مفرد: «بگو»

◇ مصدر پیشوندی «وایسادن» نیز برای ساخت امری مفرد و جمع دو بن واژه متفاوت دارد، با این تفاوت که هر دو بن واژه آن غیررسمی است. امر مفرد: «وایسا»، امر جمع: «وایسید».

◇ از بن مضارع گریستن (گری) فعل امر مفرد ساخته نمی‌شود و در حالت امر صرفاً از فعل مرکب استفاده می‌شود؛ به‌عنوان مثال: «گریه کن»

#### ۴-۲-۲. فعل‌های ناقص

برخی بن فعل‌های غیررسمی مانند فعل‌های ناقص که همه ساخت‌ها و زمان‌هایشان متداول نیست (انوری و گیوی ۱۳۹۱)، در همه ساخت‌ها استفاده نمی‌شود و در صورتی می‌تواند در ساخت فعل مورد استفاده قرار گیرد که همراه با پیشوندهای تصرفی (می، ن) باشد، در غیر این صورت، بیشتر، از بن رسمی استفاده می‌شود. فعل‌های پیشوندی که با این بن‌ها ساخته می‌شوند نیز الگویی مشابه دارند (جدول ۳).

### جدول ۳. فعل‌های ناقص ساخت‌های ماضی

بن ناقص فعل ماضی	داشت (گذاشت)	شوند <sup>۱</sup> (نشاند)	نشست (نشست)	واذاشت (واگذاشت)	فروشت (فرونشست)
ساخت مجاز ساده منفی	نداشتم	نشوندی	نشستم	وانداشت	فرونشست
استمراری	(ن)میداشتم	(ن)میشوندی	(ن)میشستم	وا(ن)میداشت	فرو(ن)میشست
نقلی مستمر	(ن)میداشته	(ن)میشونده	(ن)میشسته	وا(ن)میداشته	فرو(ن)میشسته
نقلی بعید التزامی-منفی	نداشته	نشونده	نشسته	وانداشته	فرونشسته
ساخت مجاز غیرمجاز ساده	* داشتم	* شوندی	* نشستم	* واداشتم	* فروشت
نقلی بعید التزامی-مثبت	* داشته	* شونده	* نشسته	* واداشته	* فروشته

#### ۴-۲-۳. اتصال پی‌بست‌های ضمیری فعل مختوم به واکه

«مگردومیان» تغییرات آوایی و نگارشی را هنگام اتصال پی‌بست‌های ضمیری به افعال غیررسمی مضارع اخباری و ماضی نقلی سوم شخص مفرد بررسی کرده است (Megerdooian 2008). با توجه به این که این افعال به واکه «( از لحاظ نگارشی «ه» ختم می‌شوند، در هنگام تلفظ، واکه به همخوان «ت» تبدیل می‌شود. این تغییر آوایی در صورت نگارشی، با جایگزینی «ه» با «ت» و یا درج «ت» در واژه نمایش داده می‌شود. (۸). همه بچه‌ها رو صحیح و سالم با اتوبوس برده و اوردتَشون/ اُورده\_تَشون.

#### ۵. قواعد و واژگان غیرفعلی

در این بخش به واژگان غیرفعلی تحلیلگر، شامل ستاک‌های غیرفعلی و وندهای تصریفی و پی‌بست‌های متصل به آن‌ها پرداخته و قواعد تصریفی آن بررسی می‌شود.

#### ۵-۱. واژگان

واژگان در این بخش شامل ستاک‌های غیرفعلی، وندهای تصریفی و پی‌بست‌های رسمی و غیررسمی است که می‌توانند به هر ستاک متصل شوند. ستاک‌ها بر اساس دستور زبان فارسی «انوری و گیوی» (۱۳۹۱) دسته‌بندی و مقولات واژگانی اسم، قید،

۱. با توجه به داده‌های پیکره برای بن «نشاند» دو تک‌واژه گونه غیررسمی «نشوند» و «شوند» شناسایی و به تحلیلگر معرفی شده است.

عدد، حرف ربط، حرف اضافه، شاخص، ضمیر (مشترک، شخصی)، ضمیر / صفت (اشاره، مبهم، پرسشی، تعجبی) متمایز شده‌اند. در بخش ضمیر / صفت در صورتی که ستاک به تنهایی استفاده شود، صفت و یا ضمیر است (مستقل از بافت قابل تشخیص نیست) و در غیر این صورت، ضمیر است. این نکته نیز قابل ذکر است که برچسب‌دهی به مقولات واژگانی مستقل از بافت تعریف شده است. جدول ۵، تعداد ستاک‌هایی را که ذیل هر مقوله واژگانی برای تحلیلگر تعریف شده، نشان می‌دهد.

مانند بخش افعال، در این بخش نیز در صورتی که واژه ساختار یا حداقل یک عنصر غیررسمی داشته باشد، به‌عنوان واژه غیررسمی شناسایی و پردازش می‌شود؛ به‌عنوان مثال: «کتابا، جعبه‌مون». همچنین، برای تمام اسامی، اطلاعاتی در مورد مفرد یا جمع بودن هر اسم، به‌عنوان مثال اسامی جمع مکسر عربی و نیز وندهای جمعی که هر اسم می‌تواند بپذیرد، مثلاً «ان، ها» برای اسامی جاندار، وارد شده است.

جدول ۴. مقولات واژگانی غیرفعلی تعریف شده برای تحلیلگر تصریفی

مقوله واژگانی	اسم	اسم خاص	قید	صفت	عدد
تعداد رسمی	۲۷۹۴۰	۶۲۴۲	۱۳۷۴	۱۶۶۵۰	۴۸
غیررسمی	۸۰	۱۸	۲۲۷	۴۰	۱۷
مقوله واژگانی	صفت ضمیر	شاخص	ضمیر	حرف ربط	حرف اضافه
تعداد رسمی	۷۷	۲۹	۱۰	۳۰۲	۱۶۴
غیررسمی	۲۲	۳	۳	۲۶	۴

#### ۵-۱-۱. تک‌واژه‌های جمع

در زبان فارسی اسم‌ها با وندهای تصریفی «ها»، «ان» و «گان (تک‌واژه‌گونه ان)» جمع بسته می‌شود. افزون بر آن، وندهای جمع عربی مانند «ین، ون، ات (جات)» در زبان فارسی رایج و مورد استفاده است. در ساختار غیررسمی نیز افزون‌بر وندهای فوق، از «ا» به‌تنهایی استفاده می‌شود.

(۹). سلام دوستان! طرَحای جدیدمون آماده‌س.

۱-۲. تک‌واژه‌های نکره/موصول، کسره اضافه و معرفه‌ساز

جدول ۵. الگوهای نگارشی تک‌واژه‌های نکره/موصول، کسره اضافه و معرفه‌ساز

سناک	مختوم به همخوان	مختوم به واکه الف/واو	مختوم به واکه ه	مختوم به واکه ی
تک‌واژه کسره اضافه -	ای	Ø / ی	برده، برده‌ی	-
تک‌واژه نکره	ای	ایی / ائی / اعی	ای / آئی / آعی	ای / آعی / آعی
	کتابی	کتابایی، کدوئی، طلاعی، (یه) جوراعی، ایتالیاعی	پرده‌ای، خنده‌ئی، گره‌بی	دوگانگی‌ای، اپیدمی‌ی، آگاهی‌عی
تک‌واژه معرفه‌ساز	اه	ائه / اعه	اه / آهه / آئه / آعه	اه
	کتابه	آب‌نمائه (عه)، تابلوئه (عه)	گیره‌ه، پنبه‌هه، گشته‌ئه، قراضه‌ئه	نظامیه

تک‌واژه معرفه‌ساز «ه» با اتصال به پایه اسم یا صفت جایگزین اسم، آن را شناسه می‌کند و منحصر به گونه غیررسمی است. جدول ۶، الگوهای نگارشی را برای تک‌واژه‌های نکره/موصول، کسره اضافه و معرفه‌ساز در گونه غیررسمی، با توجه به بافت حرفی که در آن قرار می‌گیرند، نشان می‌دهد.

۱-۳. پی‌بست‌های ربطی

این پی‌بست ساخت اسنادی را به اسم‌ها و غالب اقسام دیگر واژه که در این بخش آمده، اضافه می‌کند. جدول ۶، الگوهای مختلف نگارشی را در بافت حرفی نشان می‌دهد.



### جدول ۶. پی‌بست‌های ربطی<sup>۱</sup>

ش.ش.	رسمی بودن	همخوان	واکه واو	واکه ه	واکه الف	واکه ی
مفرد ۱، ۲، ۳	رسمی	س! (م، ی، ه)	س! ای (م، ی)	س! (م، ی، ه)	س! ای (م، ی)	س! (م، ی، ه)
		ست	س! است	ست	س! است	س! است
	غیررسمی	س! (م، ی، ه)	س! ای (م، ی)	س! (م، ی، ه)	س! ای (م، ی)	س! (م، ی، ه)
		مریضم، مریضی، مریضه	س! ای س! نه بدگوم، بدگوم، بدگوی، بدگوته،	س! (م، ی، ه) س! (م، ی، ه) برده م، برده یی، برده س	س! ای (م، ی) س! ای س! اس هیولام، بالام، بالای، بالاس	س! (م، ی، ه) س! (م، ی، ه) س! است نظامی م، نظامی ای، نظامی ی نظامیست
جمع ۱، ۲، ۳	رسمی	س! (یم، ید، ند)	س! ای (یم، ید، ند)	س! (یم، ید، ند)	س! ای (یم، ید، ند)	س! (یم، ید، ند)
		ند	ند	اند	ند	ند
	غیررسمی	س! (یم، ین، ن)	س! ای (یم، ین، ن)	س! (یم، ین، ن)	س! ای (یم، ین، ن)	س! (یم، ین، ن)
		مریضم، مریضین، مریضن	س! ای (یم، ین، ن) س! ای (یم، ین، ن) پرروایم، بدگویاید، بدگویین، پررویین، پررواند،	س! (یم، ین، ن) س! (یم، ین، ن) ان، ن برده یم، برده یید، برده ین، برده یین، پررون، پررون، برده ن	س! ای (یم، ین، ن) س! ای (یم، ین، ن) تنها یم، تنها یاید، تنها یین، تنها اند، تنهان	س! (یم، ین، ن) س! (یم، ین، ن) س! (یم، ین، ن) نظامی یم، نظامی یاید، نظامی یین، نظامی ان، نظامین

### ۵-۴. پی‌بست‌های ضمیری (شخصی)

این پی‌بست‌ها می‌توانند در نقش ملکی یا مفعولی ظاهر شوند و نیز با اتصال به جزء غیرفعلی فعل مرکب نقش مطابقت مفعولی یا فاعلی را به عهده گیرند (شقایق ۱۳۹۴). از آنجا که این تمایز بیشتر در سطح واژه و تصریف قابل تمیز نیست، همه آن‌ها در بخش واژه‌های غیرفعلی با عنوان پی‌بست‌های ضمیری (شخصی) نام‌گذاری شده‌اند.

۱. الگوهای نگارشی متفاوت برای یک تک‌واژ داخل گروه آورده شده است.

جدول ۲. پی‌بست‌های ضمیری

شخص شمار رسمی بودن	سناک مختوم به همخوان	سناک مختوم به واکه او- واکه الف	سناک مختوم به واکه ه	سناک مختوم به واکه ی
رسمی	س! (م،ت،ش)	س! ی (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)
غ.رسمی	کتا ب م	س! (م،ت،ش)	س! (م،ت،ش)	سبزی م
رسمی	س! (مان،تان،شان)	س! ی (مان،تان،شان)	س! ی (مان،تان،شان)	س! (مان،تان،شان)
غ.رسمی	س! (مون،تون،شون)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مون،تون،شون)
	کتا بشون	س! (مون،تون،شون)	س! (مون،تون،شون)	سبزی شون
		لیوشان، لیوتون، کتیراتان، کتیراشون	جعبه تان، جعبه شون	

۵-۱-۵. سایر تک‌واژه‌های وابسته مختص گونه غیررسمی

«شقاقی» پی‌بست‌های ساده زبان فارسی را که صورت کوتاه‌شده تک‌واژه‌های مستقل «را، و، هم» است و بیشتر در محاوره و گفتار سریع و خودمانی به کار می‌رود، معرفی می‌کند (۱۳۹۴). این پی‌بست‌ها در گونه نوشتاری غیررسمی نیز مطابق جدول ۸، به‌عنوان آخرین تک‌واژه به ساختار واژه افزوده می‌شود. افزون بر موارد ذکر شده، با توجه به بررسی داده‌ها در گونه غیررسمی، پی‌بست‌های تأکیدی «ها، ا» که «شقاقی» تنها قابل اتصال به گروه فعلی می‌داند، می‌تواند به مقوله‌های غیرفعلی نیز افزوده شوند که در واژگان غیرفعلی تحلیل‌گر تعریف شده است.

(۱۰). آدم خیلی تنبلی بود، تنبلا!، باورت همیشه!، دست به سیاه سفید نمی‌زد.

جدول ۸. برخی تک‌واژه‌های وابسته مختص گونه غیررسمی<sup>۱</sup>

سناک مختوم به ...	ها (تاکید)	ا (تاکید)	و (عطف)	رو (را)	و (را)	رم (ر+هم)	م (هم)
س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)	س! (م،ت،ش)
کتا ب م	کتا ب م	کتا ب م	کتا ب م	کتا ب م	کتا ب م	کتا ب م	کتا ب م
س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)	س! (مان،تان،شان)
س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)	س! (مون،تون،شون)
کتا بشون	کتا بشون	کتا بشون	کتا بشون	کتا بشون	کتا بشون	کتا بشون	کتا بشون
لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون	لیوشان، لیوتون، کتیراتان، کتیراشون

۱. س: سناک

مختوم به ...	ستاک	ها (تاکید)	۱ (تاکید)	و (عطف)	دو (را)	و (را)	رم (را+هم)	م (هم)
واکه الف	س!ها کتابها	-	س!او کتاباو	س!او کتاباو	س!ارو کتابارو	س!او کتاباو	س!ارم کتابارم	س!م کتابام
واکه ه	س×ها پنبه‌ها	-	س×او پنبه‌او	س×او پنبه‌او	س×ارو پنبه‌ارو	س×او پنبه‌او	س×ارم پنبه‌ارم	-

## ۲-۵. قواعد غیر فعلی

ساختار تصرفی مقوله‌های واژگانی یا به عبارتی نوع و ترتیب پیوستن وندهای تصرفی و پی‌بست‌ها به ستاک‌های غیر فعلی در گونه رسمی زبان با توجه به پژوهش‌های پیشین، از جمله «اسلامی» و همکاران (۱۳۸۳) و «انوری و گیوی» (۱۳۹۱) تعیین شده و در تحلیلگر مورد استفاده قرار گرفته است. برای گونه غیر رسمی نیز تلاش شد با توجه به پیکره جمع‌آوری شده و سایر پژوهش‌های این حوزه از جمله Megerdooian (2008) برای هر یک از مقوله‌های واژگانی قواعد جامع و مانع به گونه‌ای تعریف شود که هم تمام ساخت‌های ممکن برای هر مقوله را شامل شود و هم از تولید ساخت‌های نادرست و زاید جلوگیری کند. این قواعد به‌طور خلاصه در جدول‌های ۹ و ۱۰ آورده شده است. جدول ۹، ترتیب و امکان باهم آیی را برای وندها و پی‌بست‌های گوناگون نشان می‌دهد که در ۹ گروه تعریف شده است. جدول ۱۰، ساخت‌های ممکن برای هر مقوله واژگانی را با استفاده از گروه‌های تک‌واژی معرفی شده در جدول ۹، توصیف می‌کند. در هر دو جدول تک‌واژ یا تک‌واژه‌هایی که داخل پرانتز قرار گرفته‌اند، اختیاری‌اند و امکان حذف از ساختار واژه را دارند. در جدول ۹، تک‌واژ یا تک‌واژه‌هایی که نسبت به هم در توزیع تکمیلی هستند، در ردیف‌های مجزا، تعریف شده و در جدول ۱۰ با نشانه «/» از هم جدا شده‌اند. به‌عنوان مثال، جدول ۱۰ در مقوله «حرف اضافه» نشان می‌دهد که پی‌بست ضمیری می‌تواند به حرف اضافه متصل شود و بعد از آن پی‌بست ربطی و تأکیدی (ساخت ۲) قرار گیرد؛ «حواست بهشونه‌ها!». همچنین، این امکان وجود دارد که بعد از پی‌بست ضمیری، پی‌بست کوتاه شده «هم:م» قرار گیرد که با پی‌بست ربطی و تأکیدی در توزیع تکمیلی است؛ «حواست بهشونم (بهشون هم) باشه». همین‌طور در «تو خودت جیره‌خورشونی!» اسم بعد از همراه شدن با ساختار پی‌بست ضمیری (شون) با یک پی‌بست ربطی دوم شخص مفرد (ی) همراه می‌شود (جدول ۹، شماره ۹).

جدول ۹. الگوهای باهم آیی وندهای تصریفی و بی بست ها در گونه غیر رسمی

شماره	باهم آیی تک واژه ها	شماره	باهم آیی تک واژه ها
۱	را + (تاکید)	۲	(ربطی) + (تاکید)
۳	(را) + (هم) *	۴	نکره + (ربطی)
۵	را + (هم * تاکید)	۶	ضمیری + (ربطی را هم)
۷	(تکواژ عالی ساز تکواژ تفضیلی ساز)	۸	موصولی + (ربطی را + (هم) ... + که
۹	(جمع) + موصولی + (ضمیری + (ربطی هم) ... + که (ربطی) + (تاکید) را + (تاکید) هم عطف اضافه (ربطی) + (تاکید) را + (تاکید) هم عطف اضافه (هم * تاکید) نکره + (ربطی) + (عطف تاکید) ضمیری + (ربطی) + (تاکید) هم عطف اضافه (هم * تاکید) + (ربطی) + (تاکید) را + (تاکید) هم عطف اضافه معرفه + (را) + (هم * تاکید)		

## جدول ۱۰. الگوهای باهم‌آیی و نندهای تصریفی و پی‌بست‌ها در ساختار مقوله‌های واژگانی غیرفعلی در گونه غیررسمی

مقوله واژگانی ساختار	مقوله واژگانی ساختار	مقوله واژگانی ساختار	مقوله واژگانی ساختار
شخصی	شخصی + (۸ / ۵ / ۲ / عطف / اضافه)	اسم	اسم + ۹
پرشش	پرشش + جمع + (۶ / ۴ / ۳ / ۲)	صفت	صفت + ۷ + ۹
میهم	میهم + ضمیری + (۲ / ۱ / هم)	عدد	عدد + تا + (جمع)
میهم + (۵ / ۲ / عطف / اضافه)	عدد + (ترتیبی ۱) + (ترتیبی ۲) + (جمع)		
مشترک	مشترک + (ضمیری) + (۲ / ۱ / هم / عطف) مصدر	(پیشوند) + (منفی) + بن ماضی + ۹	
حرف اضافه	حرف اضافه + (ضمیری) + (۲ / هم)	صفت مفعولی (پیشوند) + (منفی) + صفت مفعولی + ۹	
اشاره <sup>۱</sup>	اشاره + ۹		

## ۶. ارزیابی

به‌منظور ارزیابی تحلیلگر تصریفی فارسی معاصر، از حدود ۳۰ درصد پیکره فارسی معاصر که در بخش ۳، توضیح داده شد، استفاده شده است. به این منظور، جملاتی از هر سیاق زبان فارسی غیررسمی به‌صورت تصادفی به‌گونه‌ای انتخاب شد که نسبت تنوع گونه و سیاق در آن‌ها حفظ شود و خطای سوگیری و خطای تصادفی کنترل‌شده به حداقل برسد (Biber 1993a/1993b). با توجه به این که تحلیل واژه‌ها در تحلیلگر، مستقل از بافت انجام می‌شود، نخست، واژه‌های جملات به‌صورت دستی جداسازی و استخراج شده، ۱۷۸۶ واژه یکتای به‌دست‌آمده برای آزمون تحلیلگر مورد استفاده قرار گرفت. در ارزیابی رویکرد مستقل از بافت، توانایی تحلیلگر در تولید درست تمام قواعد ممکن برای هر واژه سنجیده می‌شود. هر یک از این قواعد تعریف‌کننده تایپ<sup>۲</sup> خاصی از واژه است که در بافت مخصوص به خود معنا پیدا کرده و استفاده می‌شود. هر تایپ قاعده‌سازنده مختص به خود دارد. به‌طور مثال، برجسب‌هایی که در پیکره متنی زبان فارسی (Bijankhan et. Al. 2011) استفاده شده، مختص به قاعده‌سازنده (تایپ) آن کلمه در بافتی خاص است، اما قاعده(های) سازنده‌ای که تحلیلگر در خروجی خود تولید می‌کند، مربوط به تمام بافت‌هایی است که کلمه می‌تواند در آن‌ها ظاهر شده و به نسبت آن

۱. اشاره شامل ساختار (کسره) اضافه و معرفه نمی‌شود.

2. type

معانی مختلف داشته باشد. به عبارت دیگر، تمام تایپ‌های کلمه پوشش داده می‌شود. جدول ۱۱، خروجی تحلیلگر را برای ورودی «مردم» نشان می‌دهد. تمام تایپ‌های ممکن این ورودی در بافت‌های مختلف در میان خروجی آن قرار دارد. از میان خروجی‌ها تایپ اول و سوم غیررسمی و بقیه رسمی‌اند.

جدول ۱۱. نمونه خروجی تحلیلگر تصریفی برای ورودی «مردم»

برجسب اجزای سخن تایپ	مثال تایپ در بافت	ورودی خروجی
فعل ماضی ساده	اگه مردم (مرد هم) جوابشو نده.	<ف.م.س=مرد+ش+۳+هم>
فعل ماضی ساده	اگه مردم (من مردم) سر خاکم بیا.	<ف.م.س=مرد+ش+۱+رسمی>
اسم عام	مردم (مرد هم) مردای قدیم.	<اسمعام=مرد+هم>
اسم عام	مردم (مرد من) غیر تیه.	<اسمعام=مرد+وشخصی+۱+رسمی>
اسم عام	صدام نشون میده که مردم (مرد هستم).	<اسمعام=مرد+وربطی+۱+رسمی>
اسم عام	مردم جمع شده بودن.	<اسمعام=مرد+رسمی>

در این رویکرد، مثبت‌های درست<sup>۱</sup> قاعده‌هایی هستند که توسط تحلیلگر به درستی تولید شده‌اند. مثبت‌های اشتباه<sup>۲</sup> قاعده‌هایی هستند که به اشتباه تولید شده و منفی اشتباه<sup>۳</sup> قاعده‌هایی هستند که می‌بایست توسط تحلیلگر تولید می‌شده، اما تولید نشده‌اند. فراخوانی<sup>۴</sup> نشان‌دهنده میزان قاعده‌های درست شناخته شده برای قطعات است و از تقسیم TP بر جمع TP و FN به دست می‌آید. صحت<sup>۵</sup> نشان‌دهنده میزان درست بودن قواعدی است که توسط تحلیلگر تولید شده‌اند و از تقسیم TP بر جمع TP و FP به دست می‌آید. معیار اف<sup>۶</sup> نیز از ترکیب دو معیار فراخوانی و صحت به شکل متوازن به دست می‌آید، به طوری که می‌توان میزان تأثیر هر یک از دو معیار را مشخص کرد. رابطه (۱) توزیع مساوی (۵۰ درصد- ۵۰ درصد) هر دو معیار را به دست می‌دهد.

$$\text{معیار اف} = 2 \times \frac{\text{درستی} \times \text{فراخوانی}}{\text{درستی} + \text{فراخوانی}} \quad (1)$$

جدول ۱۲، نتایج ارزیابی تحلیلگر را با استفاده از سه معیار فراخوانی، صحت و معیار اف برای واژه‌های رسمی و غیررسمی آزمون نشان می‌دهد.

- |                       |                        |                        |
|-----------------------|------------------------|------------------------|
| 1. true positive (TP) | 2. false positive (FP) | 3. false negative (FN) |
| 4. recall             | 5. precision           | 6. f-measure           |

### جدول ۱۲. نتایج ارزیابی تحلیگر برای گونه رسمی، غیررسمی و مجموع هر دو گونه

گونه زبانی	فراخوانی	صحت	معیار اف
رسمی	۹۵/۱	۹۹/۴۱	۹۷/۲۱
غیررسمی	۹۶/۲۳	۹۹/۱۶	۹۷/۶۷
معاصر (مجموع)	۹۵/۵۶	۹۹/۶۵	۹۷/۵۶

تحلیگر تصریفی برای مجموع ۱۷۸۶ واژه یکتای آزمون، ۳۷۰۴ قاعده یا تایپ را با موفقیت (TP) تولید کرده و همین‌طور، ۱۷۲ قاعده را می‌بایست تولید می‌کرده، اما از دست داده (FN) است و ۱۳ قاعده را هم به‌اشتباه (FP) تولید کرده است.

جدول ۱۳، بررسی موردی خطای FN را به تفکیک برای دو گونه رسمی و غیررسمی نشان می‌دهد. همان‌طور که مشاهده می‌شود، کمبود واژگان با وزن ۷۱/۵ درصد، مؤثرترین عامل در رخداد خطای FN در بین همه تایپ‌هاست. رتبه‌های بعدی شامل خطای املائی (۲۵ مورد، ۱۴/۵ درصد کل عوامل) و نقص در قاعده (۱۰ مورد، ۵/۸ درصد) است.

### جدول ۱۳. بررسی خطای FN برای هر دو گونه زبانی رسمی و غیررسمی

گونه زبانی	کمبود واژگان	نقص در قاعده	خطای هم‌صدا و کسره اضافه	به‌هم‌چسبیده	تغییر آوایی	خطای املائی
رسمی	۸۳	۶	۰	۱	۰	۱۷
غیررسمی	۴۰	۴	۳	۳	۵	۸

به‌علت نمونه‌گیری تصادفی در جمع‌آوری پیکره و داده آزمون، نسبت واژه‌های رسمی آزمون به واژه‌های غیررسمی در حدود ۶۰ به ۴۰ است. از طرف دیگر، دایره واژه‌های رسمی گسترده‌تر از واژه‌های غیررسمی است و این منجر به تولید واژه‌های خارج از واژگان بیشتر رسمی نسبت به گونه غیررسمی می‌شود (طبق جدول ۱۳، دو برابر؛ ۸۳ به ۴۰). واژه‌های خارج از واژگان فراخوانی تایپ‌های رسمی را بیش از غیررسمی کاهش داده است. از طرف دیگر، ساخت‌های ساده‌تر رسمی صحت تایپ‌های رسمی را نسبت به غیررسمی افزایش داده است.

در این پژوهش جداسازی واژه‌های آزمون به‌صورت دستی انجام شده است. پیش‌بینی می‌شود اگر این جداسازی به شکل خودکار انجام شود، فراخوانی را به‌ویژه برای واژه‌های غیررسمی کاهش دهد.

## ۷. نتیجه‌گیری

تحلیلگر تصریفی پایین‌ترین سطح پردازش را (بعد از جداساز) بر روی متن فارسی انجام می‌دهد و اطلاعات تصریفی را به قطعه‌های متن خام اضافه می‌کند. اطلاعات به‌دست‌آمده از تحلیلگر تصریفی را می‌توان در سطوح بالاتر تحلیل نحوی و شناسایی گروه‌های نحوی استفاده کرد. با توجه به این‌که واژه‌های رسمی بخش قابل توجهی از واژه‌های فارسی امروز را تشکیل می‌دهد، تحلیل تصریفی این واژه‌ها نیز در کنار واژه‌های فارسی غیررسمی بخشی از فرایند پردازشی در تحلیلگر تصریفی فارسی معاصر است. برای تحلیل تصریفی همهٔ این نوع واژگان فارسی از مبدل حالت محدود استفاده شده است. به‌دلیل رویکرد مستقل از بافت تحلیلگر، برای هر واژه تمام تحلیل‌ها و قواعد احتمالی مربوط به خوانش و تحلیل معنایی، وابسته به بافت‌های گوناگون ارائه می‌شود. با افزودن یک مدل آماری برچسب‌زن صرفی-نحوی<sup>۱</sup> می‌توان این ابزار را به یک تحلیلگر تصریفی مبتنی بر بافت تبدیل کرد. در مرحلهٔ شناسایی قاعدهٔ دقیق از بین قاعده‌های تولیدی برای هر واژه، جدا از روش‌های آماری و مبتنی بر بافت، می‌توان از اطلاعات افزوده‌ای که در تصریف غیررسمی نسبت به تصریف رسمی وجود دارد، کمک گرفت. وندها و واژه‌بست‌هایی که امکان اتصال به ستاک در ساختار غیررسمی را دارند، به‌دلیل تنوع بالاتر در بردارندهٔ اطلاعات صرفی-نحوی بیشتری هستند که می‌تواند در تجزیه و تحلیل نحوی و معنایی مورد استفاده قرار گیرد. همچنین، خروجی این تحلیلگر به‌گونه‌ای طراحی شده است که بتوان ابزارهای مکمل و ساده‌ای مانند ریشه‌یاب و ابزار تقطیع ساخت‌واژی بر اساس آن ساخت تا در تحلیل زبانی به کار آید. از جمله کاربردهای دیگر این ابزار می‌توان به آموزش فارسی به غیرفارسی‌زبانان اشاره کرد. با توجه به پوشش فارسی غیررسمی، آموزش ساختار تصریفی کلمات گفتار به زبان‌آموزان فارسی و تمرین و تکرار آن با این ابزار توسط زبان‌آموز امکان‌پذیر است.

بر اساس جدول ۱۲، بیشترین منفی‌های اشتباه در فارسی رسمی و در فارسی غیررسمی مربوط به کمبود واژگان است. با افزودن واژگان مناسب تا اندازهٔ قابل توجهی فراخوانی تحلیلگر افزایش خواهد یافت. بعد از نقص کمبود واژگان بیشترین تعداد منفی اشتباه مربوط به خطاهای املائی است. در صورت استفاده از یک خطیاب و اصلاح‌کنندهٔ خطای املائی در کنار تحلیلگر تصریفی، فراخوانی به‌ویژه برای گونهٔ رسمی افزایش خواهد یافت.

1. morpho-syntactic part-of-speech tagger



## فهرست منابع

- اسلامی، محرم، مسعود آتشگاه، صدیقه علی‌زاده لمجیری، و طاهره زندی. ۱۳۸۳. واژگان زبانی فارسی. اولین کارگاه پژوهشی زبان فارسی و رایانه. تهران.
- اسلامی، محرم، و صدیقه علی‌زاده لمجیری. ۱۳۸۸. ساختار تصرفی کلمه در زبان فارسی. زبان و ادب فارسی. نشریه دانشکده ادبیات و علوم انسانی دانشگاه تبریز ۲۱۱: ۱-۱۸.
- انوری، حسن، و حسن گیوی. ۱۳۹۱. دستور زبان فارسی. ویرایش چهارم. تهران: انتشارات فاطمی.
- بی‌جن‌خان، محمود. ۱۳۸۶. پیاده‌سازی استاندارد اینگلس در پیکره متنی زبان فارسی معاصر. مطالعه و تحقیق جهت تدوین پژوهشنامه عملیاتی دادگان. تهران: دبیرخانه شورای عالی اطلاع‌رسانی.
- تازه‌جانی، سمیه، و محمد بحرانی. ۱۳۹۲. بررسی روند تغییرات در تبدیل فرم رسمی افعال به فرم محاوره‌ای و ارائه یک تحلیلگر صرفی برای افعال محاوره‌ای. اولین هم‌اندیشی زبان فارسی و اینترنت. تهران.
- حیدرپور، داود، مصطفی مصالحی، محمود بی‌جن‌خان، هادی ویسی، و وحید رنجبر. ۱۳۹۸. شناسایی و پوشش واحدهای خارج از واژگان در فارسی غیررسمی. پنجمین همایش ملی زبان‌شناسی رایانشی. انجمن زبان‌شناسی ایران، تهران.
- شقایق، ویدا. ۱۳۹۴. نندگروهی. زبان و زبان‌شناسی. مجله انجمن زبان‌شناسی ایران، پژوهشگاه علوم انسانی و مطالعات فرهنگی ۹ (۱۷): ۱-۲۶.
- صادقی، علی‌اشرف، و زهرا زندی‌مقدم. ۱۳۸۵. فرهنگ املائی زبان فارسی. تهران: فرهنگستان زبان و ادب فارسی.
- طیب‌زاده، امید. ۱۳۹۸. مبانی و دستور خط فارسی شکسته. تهران: پژوهشگاه علوم انسانی و مطالعات فرهنگی.
- کاشفی، امید، میترا نصری، و کامیار کنعانی. ۱۳۸۹. خط‌یابی املائی خودکار در زبان فارسی. تهران: دبیرخانه شورای عالی اطلاع‌رسانی.
- مواجی، وحید، محرم اسلامی، و بهرام وزیرنژاد. ۱۳۹۰. پارس مورف: تحلیلگر صرفی زبان فارسی. پردازش علائم و داده‌ها ۸ (۱): ۳-۸.

## References

- Alshargi, Faisal, and Owen Rambow. 2016. Morphologically Annotated Corpora and Morphological Analyzers for Moroccan and Sanaani Yemeni Arabic. In 10th Language Resources and Evaluation Conference (LREC 2016). Portoroz.
- Badawi, El Said, Michael Carter, and Adrian Gully. 2013. *Modern Written Arabic: A Comprehensive Grammar*. London: Routledge.
- Biber, Douglas. 1993a. Representativeness in Corpus Design. *Literary and Linguistic Computing* 8 (4): 57-243.
- Biber, Douglas. 1993b. Using Register-Diversified Corpora for General Language Studies. *Computational Linguistics* 19 (2): 41-219.

- Bijankhan, Mahmood, Javad Sheykhzadegan, Mohammad Bahrani, and Masood Ghayoomi. 2011. Lessons from building a Persian written corpus: Peykare. *Language Resources and Evaluation* 45 (2): 143–164.
- Ferguson, Charles A. 1959. Diglossia. *Word* 15 (2): 325–340.
- Hulden, Mans. 2009. Foma: A Finite-State Compiler and Library. In *Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics: Demonstrations Session*: 29–32. Association for Computational Linguistics. Athens.
- Jurafsky, Daniel, and James H Martin. 2008. *Speech and Language Processing*. 2nd Edition. New Jersey: Prentice Hall.
- Lazard, G. 1992. *A grammar of contemporary Persian*. Costa Mesa, CA: Mazda Publishers.
- McEnery, Tony, and Andrew Hardie. 2011. *Corpus Linguistics: Method, Theory and Practice*. Cambridge: Cambridge University Press.
- Megerdooian, Karine. 2008. *Analysis of Farsi Weblogs*. Washington DC: MITRE Corporation,.
- \_\_\_\_\_. 2000. *Persian Computational Morphology: A Unification-Based Approach*. New Mexico: Computing Research Laboratory, New Mexico State University.
- Ryding, Karin C. 2005. *A Reference Grammar of Modern Standard Arabic*. Cambridge: Cambridge university press.
- Sagot, Benoit, and Geraldine Walther. 2010. A morphological lexicon for the Persian language. In *Proceedings of the 7th Language Resources and Evaluation Conference (LREC'10)*. Malta.
- Shamsfard, M., H. S. Jafari, and M. Ilbeygi. 2010. *STeP-1: A Set of Fundamental Tools for Persian Text Processing*. In *LREC 2010-8th Language Resources and Evaluation Conference*. Malta.
- Torjmen, R. and K. Haddar. 2018. *Morphological analyzer for the Tunisian dialect*. In *International Conference on Text, Speech, and Dialogue* (pp. 180-187). Springer, Cham. Czech Republic.

#### داود حیدریور

متولد سال ۱۳۶۲، دانش‌آموخته رشته نرم‌افزار و مترجمی زبان انگلیسی در مقطع کارشناسی و رشته زبان‌شناسی رایانشی در مقطع کارشناسی ارشد از دانشگاه تهران است. از جمله موضوعات تحقیقی مورد علاقه وی پردازش زبان طبیعی، متن‌کاوی و نظر‌کاوی است.



#### الهام‌سادات سیط

دارای مدرک کارشناسی فیزیک مهندسی و کارشناسی ارشد زبان‌شناسی رایانشی از دانشگاه تهران است. از جمله موضوعات تحقیقی مورد علاقه وی پردازش نحوی و ساخت‌واژی، نرمال‌سازی متن، طراحی و تهیه داده‌زبانی و مدل‌های ماتریسی در تحلیل نحوی (دستور ماتریسی) است.



### محمود بی‌جن‌خان

متولد سال ۱۳۳۷، دارای مدرک تحصیلی دکتری در رشته زبان‌شناسی از دانشگاه تهران است. ایشان هم‌اکنون استاد تمام گروه زبان‌شناسی دانشگاه تهران است.

آواشناسی، واج‌شناسی، زبان‌شناسی پیکره‌ای، تولید پیکره‌های زبانی برای آموزش و ارزیابی سامانه‌های فناوری زبان از جمله علایق پژوهشی وی است.



### مصطفی صالحی

دارای مدرک دکتری در رشته مهندسی کامپیوتر از دانشگاه صنعتی شریف است. ایشان هم‌اکنون عضو هیئت علمی دانشکده علوم و فنون نوین دانشگاه تهران است.

پایش رسانه‌های اجتماعی از جمله علایق پژوهشی وی است.



### هادی ویسی

دارای مدرک دکتری مهندسی کامپیوتر از دانشگاه صنعتی شریف است. وی هم‌اکنون عضو هیئت علمی دانشکده علوم و فنون نوین دانشگاه تهران است.

هوش مصنوعی و زبان‌شناسی رایانشی از جمله علایق پژوهشی وی است.

