



HAL
open science

Facilitating phenotyping from clinical texts: the medkit library

Antoine Neuraz, Ghislain Vaillant, Camila Arias, Olivier Birot, Kim-Tam Huynh, Thibaut Fabacher, Alice Rogier, Nicolas Garcelon, Ivan Lerner, Bastien Rance, et al.

► To cite this version:

Antoine Neuraz, Ghislain Vaillant, Camila Arias, Olivier Birot, Kim-Tam Huynh, et al.. Facilitating phenotyping from clinical texts: the medkit library. *Bioinformatics*, 2024, pp.btac681. 10.1093/bioinformatics/btac681 . hal-04802917

HAL Id: hal-04802917

<https://inria.hal.science/hal-04802917v1>

Submitted on 25 Nov 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Data and text mining

Facilitating phenotyping from clinical texts: the medkit library

Antoine Neuraz^{1,2,3}, Ghislain Vaillant^{1,2}, Camila Arias^{1,2}, Olivier Birot^{1,2}, Kim-Tam Huynh^{1,2}, Thibaut Fabacher^{1,2,4}, Alice Rogier^{1,2,5}, Nicolas Garcelon ^{1,2,6}, Ivan Lerner^{1,2,5}, Bastien Rance^{1,2,5}, Adrien Coulet ^{1,2,*}

¹Inria Paris, Paris 75013, France

²Centre de Recherche des Cordeliers, Inserm UMR 1138, Université Paris Cité, Sorbonne Université, Paris 75006, France

³Hôpital Necker, Assistance Publique—Hôpitaux de Paris, Paris 75015, France

⁴University Hospital of Strasbourg, Strasbourg 67000, France

⁵Hôpital Européen Georges Pompidou, Assistance Publique—Hôpitaux de Paris, Paris 75015, France

⁶Imagine Institute, Inserm UMR 1163, Université Paris Cité, Paris 75015, France

*Corresponding author. HeKA Team, ParisSanté Campus, 2-10 rue d'Oradour sur Glane, Paris 75015, France. E-mail: adrien.coulet@inria.fr

Associate Editor: Xin Gao

Abstract

Summary: Phenotyping consists in applying algorithms to identify individuals associated with a specific, potentially complex, trait or condition, typically out of a collection of Electronic Health Records (EHRs). Because a lot of the clinical information of EHRs are lying in texts, phenotyping from text takes an important role in studies that rely on the secondary use of EHRs. However, the heterogeneity and highly specialized aspect of both the content and form of clinical texts makes this task particularly tedious, and is the source of time and cost constraints in observational studies.

To facilitate the development, evaluation and reproducibility of phenotyping pipelines, we developed an open-source Python library named `medkit`. It enables composing data processing pipelines made of easy-to-reuse software bricks, named `medkit` operations. In addition to the core of the library, we share the operations and pipelines we already developed and invite the phenotyping community for their reuse and enrichment.

Availability and implementation: `medkit` is available at <https://github.com/medkit-lib/medkit>.

1 Introduction

The collection at large scale of Electronic Health Records (EHRs) and the constitution of Clinical Data Warehouses (CDW) enable the design of clinical studies relying on a secondary use of healthcare data (Madigan *et al.* 2014). A substantial part of the necessary information to conduct these studies is available in texts, such as clinical notes, hospitalization, or exam reports (Kharrazi *et al.* 2018). For instance, tasks such as the inclusion/exclusion of patients, and the extraction of outcome variables or covariates often require the consideration of clinical texts.

In biomedical data sciences, the two complementary tasks of either identifying individuals associated with a specific, potentially complex, trait or condition, or listing the traits of an individual are generally named *phenotyping*. And the specific case of phenotyping from clinical text is a continuous challenge for several reasons (Banda *et al.* 2018). First clinical text is highly specialized as it includes various factors of complexity such as medical entities absent from the general domain, hypotheses, negations, abbreviations, personal information; what motivates the development of dedicated phenotyping tools (Kreimeyer *et al.* 2017). Besides, many

powerful Natural Language Processing (NLP) tools and models are developed and shared for both the general and biomedical texts, making reuse, adaptation, and chaining rational approaches in biomedicine. But the highly heterogeneous aspect of clinical texts (e.g. physician versus nurse notes, hospital A versus hospital B notes, French versus English notes) makes the performance of a tool hardly predictable on a new corpus. In addition, clinical texts can barely be shared because of their personal and sensitive aspects. This implies the need for tools that ease the evaluation and adaptation of phenotyping approaches to the various types of texts generated in the large variety of existing clinical settings.

We present here `medkit`, an open-source Python library, that aims primarily at facilitating the reuse, evaluation, adaptation, and chaining of NLP tools for the development of reproducible phenotyping pipelines. By extension, `medkit` enables the extraction of information related to patient care, such as treatments or procedures, which are not phenotype *per se*. The rest of this manuscript presents the core elements of the library, develops on its easiness of use and details example pipelines developed with `medkit` for the extraction of drug treatments from clinical texts. It lists other pipelines

Received: 17 July 2024; Revised: 16 October 2024; Editorial Decision: 6 November 2024; Accepted: 12 November 2024

© The Author(s) 2024. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

already developed and ready for reuse and ends on two particularity added values of the library, which are the support of nondestructive processing and provenance tracing.

2 Related work and positioning

The PheKB initiative proposes a collaborative web portal to share phenotyping algorithms in the form of semi-formal flow charts, documenting their steps and chaining (Kirby *et al.* 2016). In this manner, PheKB helps exchanging and standardizing phenotyping algorithms. However, provided representations are not shared with their computational implementations, which limits their reproducibility and comparison. In addition, algorithmic steps that rely on clinical texts are underspecified, as they usually require an adaptation to the peculiarities of local texts. The OHDSI community offers software tools such as Atlas, which proposes standard and reusable tools for the data analysis of observational studies from EHRs (Schuemie and DeFalco 2019). Those are developed for steps coming next to the information extraction, once features are structured and normalized. medkit fills this exact step, extracting and normalizing features from unstructured parts of EHRs. The MedCAT library targets this step as well but only focuses on entity recognition and normalization with the UMLS (Kraljevic *et al.* 2021). The Gate suite provides a graphical user interface which facilitates sequential application of various preprocessing and NLP tools on texts (Cunningham 2002). But Gate is mostly adapted to educational or exploratory purposes, because of its limited ability to analyze large corpora and to adapt to novel tools. NLTK (Bird *et al.* 2009), spaCy (Honnibal and Montani 2017) and FLAIR (Akbik *et al.* 2019) are Python libraries dedicated to advanced NLP development, designed for NLP engineers and researchers. Following a different way, medkit aims at being easier to start with, facilitating the reuse and chaining of simple-to-complex NLP tools, such as those developed with the previously cited libraries.

One of the main particularity of medkit is to place a strong emphasis on nondestructive operations, i.e. no information is lost when passing data from one step to another; and on a flexible tracing of data provenance. In this matter, medkit is original and found inspiration in bioinformatic workflow management systems, such as Galaxy and Snakemake (Mölder *et al.* 2021, Community 2022), which facilitate the reproducibility of bioinformatic pipelines.

3 The core components of medkit

For internal data management, medkit represents data with three simple core classes: Documents, Annotations, and Attributes. Each of these classes is associated with properties and methods to represent data and metadata of various modalities such as audio or images, even though medkit is primarily designed for text. A Document is the minimal data structure of medkit, which associates an identifier with a set of Annotations; in turn an Annotation associates an identifier, a label, and a set of Attributes; lastly Attributes associate an identifier, a label, and a value.

For data processing, medkit defines two main classes: Operations and Pipelines. Typically, an operation is taking data as an input, runs a function over these data and returns output data. For instance, an Operation can input a Document, perform Named Entity Recognition (NER) and

output a set of Annotations associated with the Document. Accordingly, an operation can be the encapsulation of a previously developed tool, or a new piece of software developed in Python using medkit classes. Converters are particular operations for input and output management, which enable the transformation from standard formats such as CSV, JSON, Brat, Docanno annotations, into medkit Documents, Annotations and Attributes, or inversely. Lastly, Pipelines enable to chain Operations within processing workflows.

We refer readers to the medkit documentation for more details on its core components (see Availability for a web link).

4 Encapsulate, chain, and reuse operations

Numerous data processing tools exist, in particular in NLP, where pretrained models are routinely shared within libraries or platforms such as spaCy or Hugging Face (Honnibal and Montani 2017, Wolf *et al.* 2019). The goal of medkit is to facilitate their reuse, evaluation, and chaining. Following are examples of available medkit operations that reuse third-party tools: the Microsoft library named Presidio for text deidentification (Mendels *et al.* 2018); a date and time matcher from the EDS-NLP lib (Wajsburt *et al.*); text translator using transformers from the Hugging Face platform. Similarly, medkit operations enable the encapsulation of spaCy modules in particular by input, output and annotation conversion functions.

In addition, we internally developed original operations for: NER; relation extraction; preprocessing; deidentification; evaluation; the pre-annotation of clinical texts to speed-up manual annotation; the detection of negation, hypothesis, and antecedents within the context of entities; the fine-tuning of preexisting models; the classification of sentence and documents; the loading of audio patient-caregiver conversations, their diarisation and transcription to text (Nun *et al.*); and others. We aim at progressively enriching the catalogue of shared tools, thanks to the continuous growth of the community of medkit users and contributors.

The development of medkit pipelines is facilitated by three main elements. First, medkit data format asks for an initial conversion, but avoids further formatting in subsequent treatments. Second, sharing ready-to-use operations and open-source example pipelines speeds up the prototyping of new ones. Third, good practices in software development such as continuous integration, rich documentation facilitate start-off.

5 Example pipelines

As an illustration, we describe two medkit pipelines in Fig. 1 for the extraction of drug treatment from clinical texts. The first, in black, aims at comparing the performances of two NER tools named Drug NER 1 and 2, which are dictionary-based and Transformer-based methods respectively. Considering that Drug NER 2 obtained the best performances, the second pipeline is designed to use only the latter to extract the mentions of drug treatments from new texts. Both pipelines share three steps of preprocessing: conversion of raw texts into medkit Documents, sentence splitting and deidentification. The first pipeline evaluates the two tools on the basis of reference annotations saved as Brat format, whereas the second pipeline annotates new documents with

drugs and produces output annotations in Brat format. A snippet of code for the medkit implementation of the second pipeline is shown in Fig. 2. The full implementation of the two pipelines is available at https://medkit-lib.org/cookbook/drug_ner_eval/.

6 Available pipelines

We implemented and share pipelines for: the phenotyping of chemotherapy toxicities, and their grades (Rogier *et al.* 2022); the phenotyping of rheumatoid arthritis in French clinical reports (Fabacher *et al.* 2023b); the phenotyping of COVID-19 and the comparisons of pipelines relying either on

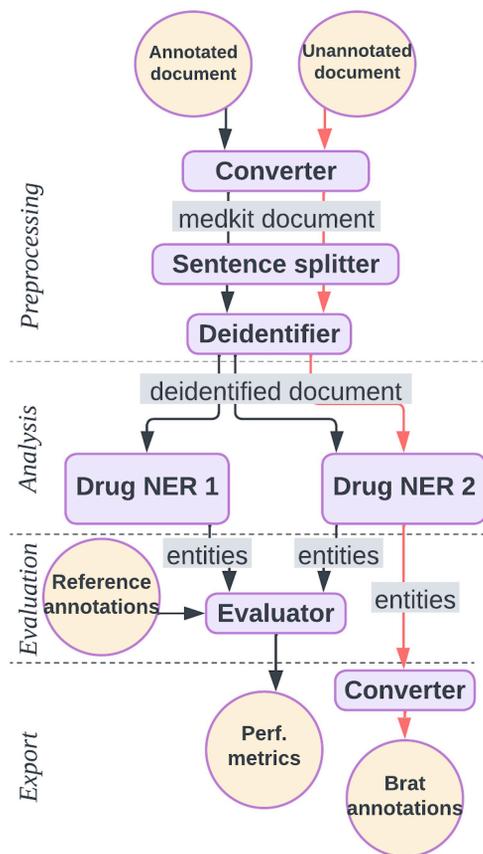


Figure 1. Example medkit pipelines. The black pipeline converts raw texts to the medkit format, deidentifies them, recognizes drug entities with two distinct tools and compute performances for comparison. The orange pipeline, performs the same preprocessing tasks, recognizes drugs with only Drug NER 2 and outputs annotations in the Brat format.

```
from medkit.core import DocPipeline
from medkit.core import Pipeline
from medkit.core import PipelineStep
from pathlib import Path
from medkit.core.text import TextDocument
#loading documents
docs = TextDocument.from_dir(path=Path("./working_dir"), pattern='*.txt', encoding='utf-8')
#pipeline definition (operations have been defined earlier. See full notebook)
my_pipeline = Pipeline(steps=[
    PipelineStep(sentence_tokenizer, input_keys=["full_text"], output_keys=["sentence"]),
    PipelineStep(deidentifier, input_keys=["sentence"], output_keys=["deided_sentence"]),
    PipelineStep(bert_matcher, input_keys=["deided_sentence"], output_keys=["drug_ner2_entities"]),
    PipelineStep(bert_matcher, input_keys=["full_text"], output_keys=["drug_ner2_entities"])
])
#execution of the pipeline on each document
doc_pipeline = DocPipeline(pipeline=my_pipeline)
doc_pipeline.run(docs)
```

Figure 2. Snippet of the implementation of the orange pipeline of Fig. 1.

the English versus French UMLS (Neuraz *et al.* 2024); the benchmarking of NER approaches on three clinical case corpora, comparing dictionary-based, transformer, and generative approaches (Hubert *et al.* 2024); the detection of text duplications in collections of clinical texts (Fabacher *et al.* 2023a). A last example of pipeline chains the recognition of drug, dates with a sentence classifier to detect treatment start and stop (Pohyer *et al.* 2024).

We refer the reader to the tutorial and cookbook sections of the medkit documentation for a list of available operations and examples of pipelines (see Availability for a web link).

7 Nondestructive processing and provenance

The medkit library features two noteworthy functionalities: nondestructive processing and flexible provenance tracing. Nondestructive processing ensures that no information is lost when passing from one operation to the next. This is of interest to get back on the raw text, after this one underwent various transformation steps, such as deidentification or character replacements. Nondestructive processing is enabled by the propagation of original spans through successive operations. We note that this functionality might be lost in the case of external and noncompliant tools encapsulated in medkit operations.

Provenance tracing consists in recording provenance data, i.e. meta-data documenting where data come from and how it was transformed. medkit implements this tracing by generating provenance data using the PROV-O standard ontology (Lebo *et al.* 2013). This tracing is flexible in the sense that users can set the level of verbosity and details they want to trace about the previous operations and states, in order to avoid generating large amounts of provenance data when those are unnecessary.

The unique combination of nondestructive processing and provenance tracing improves the explainability and reproducibility of results of pipelines of various level of complexity. These functionalities, along with its open-source nature and its focus on interoperability with existing libraries, pipelines, and models, make it well aligned with the FAIR principles for research software (Barker *et al.* 2022).

8 Availability

medkit is hosted at <https://github.com/medkit-lib/medkit>, and released under an MIT license. Its documentation, with examples and tutorials, is hosted at <https://medkit-lib.org/>.

9 Conclusion and perspectives

medkit is an open-source library for the composition of data processing pipelines made of easy-to-reuse software bricks, which aim at facilitating phenotyping from clinical texts. In addition to the core of the library, we share many of these bricks and examples of pipelines, and invite the phenotyping community for their reuse and enrichment.

So far, medkit enables linear execution of pipelines over a set of documents. Whereas it is simple to distribute the execution of pipelines by splitting a large corpus in subsets, parallelization within pipelines is not supported yet, but is planned for the future. We would like to grow the community of users of medkit, first by developing a searchable library of available operations, by enriching this library and enabling users to share their own pipelines. Pipelines may be showcased in a gallery of examples to inspire and facilitate reuse. Next developments will concern operations for the generation of features that are compliant with the OMOP Common Data Model, and operations that facilitate the use of large language models and prompting.

Acknowledgements

Authors thank users of medkit, L.-A. Guiottel, M. Hassani, S. Cossin, T. Hubert, and V. Pohyer for their insightful inputs.

Author contributions

C.A., O.B. and K.T.H.: Conceptualization, Software, Writing—review & editing. T.F. and A.R.: Software, Writing—review & editing. G.V., N.G., B.R.: Conceptualization, Project administration, Software, Supervision, Writing—review & editing. I.L.: Conceptualization, Funding acquisition, Software, Writing—review & editing. A.N. and A.C. Conceptualization, Funding acquisition, Project administration, Software, Supervision, Writing—original draft, Writing—review & editing.

Conflict of interest: No competing interest to declare.

Funding

This work was supported by Inria, Inria Paris; and the ANR under the France 2030 program [ANR-22-PESN-0007].

References

- Akbik A, Bergmann T, Blythe D *et al.* FLAIR: an easy-to-use framework for state-of-the-art NLP. In: *NAACL 2019 (Demonstrations)*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, 54–9.
- Banda JM, Seneviratne M, Hernandez T *et al.* Advances in electronic phenotyping: from rule-based definitions to machine learning models. *Annu Rev Biomed Data Sci* 2018;1:53–68. <https://doi.org/10.1146/annurev-biodatasci-080917-013315>
- Barker M, Chue Hong NP, Katz DS *et al.* Introducing the FAIR principles for research software. *Sci Data* 2022;9:622. <https://doi.org/10.1038/s41597-022-01710-x>
- Bird S, Klein E, Loper E. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. Sebastopol, California: O'Reilly Media, Inc., 2009.
- Community TG. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses. *Nucleic Acids Res* 2022;50:W345–51. <https://doi.org/10.1093/nar/gkac247>
- Cunningham H. GATE: a framework and graphical development environment for robust nlp tools and applications. In: *ACL 2002*. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, 2002, 168–75.
- Fabacher T, Birot O, Arias-Villamil C *et al.* Détection de zones dupliquées dans des comptes rendus médicaux. In: *Actes de la Journée D'étude Sur la Similarité Entre Patients*, SimPa, Paris, France, 2023a.
- Fabacher T, Sauleau E-A, Leclerc N *et al.* Evaluating the portability of rheumatoid arthritis phenotyping algorithms: case study on French EHRs. *Stud Health Technol Inform* 2023b;302:768–72. <https://doi.org/10.3233/SHTI230263>
- Honnibal M, Montani I. spaCy2: natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. 2017.
- Hubert T, Vaillant G, Birot O *et al.* Comparing NER approaches on French clinical text, with easy-to-reuse pipelines. *Stud Health Technol Inform* 2024;316:272–6. <https://doi.org/10.3233/SHTI240396>
- Kharrazi H, Anzaldi LJ, Hernandez L *et al.* The value of unstructured electronic health record data in geriatric syndrome case identification. *J Am Geriatr Soc* 2018;66:1499–507. <https://doi.org/10.1111/jgs.15411>
- Kirby JC, Speltz P, Rasmussen LV *et al.* PheKB: a catalog and workflow for creating electronic phenotype algorithms for transportability. *J Am Med Inform Assoc* 2016;23:1046–52. <https://doi.org/10.1093/jamia/ocv202>
- Kraljevic Z, Searle T, Shek A *et al.* Multi-domain clinical natural language processing with MedCAT: the medical concept annotation toolkit. *Artif Intell Med* 2021;117:102083. <https://doi.org/10.1016/j.artmed.2021.102083>
- Kreimeyer K, Foster M, Pandey A *et al.* Natural language processing systems for capturing and standardizing unstructured clinical information: a systematic review. *J Biomed Inform* 2017;73:14–29. <https://doi.org/10.1016/j.jbi.2017.07.012>
- Lebo T, Sahoo S, McGuinness D *et al.* PROV-O: the PROV ontology. *W3C* 2013;30.
- Madigan D, Stang PE, Berlin JA *et al.* A systematic statistical approach to evaluating evidence from observational studies. *Annu Rev Stat Appl* 2014;1:11–39. <https://doi.org/10.1146/annurev-statistics-022513-115645>
- Mendels O, Peled C, Vaisman Levy N *et al.* Microsoft Presidio: context aware, pluggable and customizable pii anonymization service for text and images. 2018.
- Mölder F, Jablonski KP, Letcher B *et al.* Sustainable data analysis with snakemake. *F1000Res* 2021;10:33.
- Neuraz A, Lerner I, Birot O *et al.* TAXN: translate align extract normalize, a multilingual extraction tool for clinical texts. *Stud Health Technol Inform* 2024;310:649–53. <https://doi.org/10.3233/SHTI231045>
- Nun A, Olivier B, Gaël G *et al.* Samsamu—a french medical dispatch dialog open dataset. SSRN Preprint. <https://doi.org/10.2139/ssrn.4869223>
- Pohyer V, Fabre E, Oudard S *et al.* Fake it till you predict it: data augmentation strategies to detect initiation and termination of oncology treatment. arXiv, arXiv:2410.10271, 2024, preprint: not peer reviewed.
- Rogier A, Coulet A, Rance B. Using an ontological representation of chemotherapy toxicities for guiding information extraction and integration from EHRs. *Stud Health Technol Inform* 2022;290:91–5. <https://doi.org/10.3233/SHTI220038>
- Schuemie M, DeFalco F. OHDSI analytics tools. In: *The Book of OHDSI: Observational Health Data Sciences and Informatics*, Chapter 8. OHDSI, 2019. Independently published.
- Wajsburt P, Petit-Jean T, Dura B *et al.* EDS-NLP: efficient information extraction from French clinical notes (v0.12.0). *Zenodo*, 2024. <https://doi.org/10.5281/zenodo.11238626>
- Wolf T, Debut L, Sanh V *et al.* Huggingface's transformers: state-of-the-art natural language processing. arXiv, arXiv:1910.03771, 2019, preprint: not peer reviewed.

© The Author(s) 2024. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Bioinformatics, 2024, 40, 1–4

<https://doi.org/10.1093/bioinformatics/btae681>

Applications Note