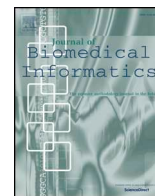




Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.



## Deep phenotyping: Embracing complexity and temporality—Towards scalability, portability, and interoperability



### 1. Introduction

Clinical data are the basic staple of health learning [1]. The rapidly growing interoperable clinical datasets, including electronic health records (EHR), administrative and claims records, and human phenotype data collected from clinical research studies, present unprecedented opportunities for developing high-throughput methods for electronic phenotyping. The term *phenotyping* is still young and evolving, emerging around 2006 when the Electronic Medical Records and Genomics (eMERGE) program was launched in the United States and directed significant effort into phenotyping using EHR data [2]. Meanwhile, the biomedical informatics research community has been exploring electronic phenotyping solutions for decades, given the essential roles electronic phenotyping plays in disease knowledge discovery, application, and clinical research. Early phenotyping research focused primarily on case ascertainment or cohort identification [3]. In contrast, deep phenotyping shifted the focus from identification to characterization, which aims to deliver precise and comprehensive characterization of observable traits representing unique morphological, biochemical, physiological, or behavioral properties of the identified patient populations [4]. Deep phenotyping brings us a step closer on the path towards Precision Medicine via the development of precise disease classification systems. Consequently, the task of deep phenotyping poses the following new requirements.

**First**, it requires the extraction of nuanced phenotypic traits, such as “short stature”, “large head”, “poor weight gain”, “depressed nasal bridge”, and “clubbing of toes”. These traits are occasionally available in structured coded data but more often can only be accurately communicated as clinical text. Therefore, natural language processing (NLP) is essential for identifying the rich information to accomplish deep phenotyping for many phenotypes [5], and has been playing an increasingly important role in deep phenotyping.

**Second**, to achieve richer, deeper, and more precise characterization, deep phenotyping algorithms need to be more expressive, semantically interoperable, and interpretable than the conventional “black box” computational solutions, which are often criticized for lacking explainability [6]. Therefore, the extracted phenotypes need to be normalized using various clinical terminologies or ontologies such as HPO (The Human Phenotype Ontology), ICD-9-CM, ICD-10-CM, and SNOMED-CT. This requirement is at odds with the first requirement because the existing semantic knowledge resources have varying concept coverage and concept granularities, posing challenges for concept normalization in addition to the limitations noted with the first requirement. Consequently, creating code sets for a phenotype has become a foundational building block and yet the most time-consuming

bottleneck in the knowledge engineering process for deep phenotyping. Since clinical databases tend to be heterogeneous, empirical knowledge of data sources and clinical processes is critical for identifying useful codes that have high sensitivity and specificity at a given site [7]. This heterogeneity also creates barriers to portability and interoperability of solutions leveraging distributed big clinical data for clinical phenotyping [8]. Standards beyond terminologies and ontologies at a higher level, such as common data models, promise to meet some of these needs, but efforts to improve existing widely adopted standards are still warranted.

**Third**, deep phenotyping requires the processing of a broad range of data types, which can include voice recordings, videos, images, genome sequences, or biological pathways, to name a few. For example, for the current COVID-19 pandemic, the ground-glass appearance in the lower lungs in chest x-rays for COVID-19 patients paired with clinical narrative reports for dry cough and dyspnea (labored breathing) are key phenotypes of COVID-19-positive patients [9–11].

**Fourth**, deep phenotyping requires more sophisticated analytics beyond case-control classification and may involve the characterization of the temporal trajectory of a phenotype or the identification of disease subtypes based on their differentiating phenotypes, progressions, and clinical outcomes. Incorporation of domain knowledge remains critical. Both supervised and unsupervised methods will be useful for deep phenotyping in different contexts.

**Fifth**, deep phenotyping requires the identification of connections between diseases and their common phenotypic traits (i.e., phenotypes that may be observed across various disease domains) to enable linking potentially common etiologies across diseases and to facilitate important applications such as drug repurposing based on disease physiology commonalities. For example, for the current COVID-19 pandemic, international doctors have reported that loss of smell or taste are widely observed in infected patients. Similarly, studies have also shown olfactory deficits are prevalent in patients with Alzheimer’s disease [12]. Therefore, the currently hot research on calculating disease similarities will play an important role in deep phenotyping to fulfill this need [13].

In order to satisfy the above requirements for advancing deep phenotyping, this special issue offers twenty original articles presenting novel methodologies for case ascertainment, patient stratification, disease subtyping and temporal phenotyping. These novel methods were demonstrated across various disease domains (such as cancer, rare diseases, obesity, acute or chronic kidney diseases, and schizophrenia, to name a few) using a broad range of novel data sources (including clinical narratives, voice, biological pathways, research questionnaire data, and claims data in addition to EHR data) while addressing

**Table 1**

Contributions of the twenty included papers in response to the five requirements for deep phenotyping.

Reference	First author	Summary of Contributions
<b>Requirement 1: Natural Language Processing</b>		
[14]	Datta, S	A systematic review of NLP on cancer notes
[15]	Liu, Q	Symptom extraction for patient stratification
[16]	Lyudovyyk, O	NLP on pathology notes for subtyping
[17]	Liu, C	Ensemble of NLP for better portability
<b>Requirement 2: Standardization</b>		
[18]	Hong, N	A FHIR-based EHR phenotyping framework
[19]	Shang, N	An empirical study of "making phenotyping work visible" that demonstrates the need for standardized processes
[20]	Hripcsak, G	Demonstrate OMOP's value in improving phenotyping algorithms' portability
[21]	Ostropolets, A	Adapting EHR phenotypes to claims data using OMOP Common Data Model
[22]	Reps, J	OMOP CDM-based probabilistic phenotyping algorithms using self-reported data
[23]	Swerdel, J	OMOP CDM-based standardized phenotype evaluation algorithms
[24]	Warner, J	Expansion of OMOP CDM to cancer phenotypes
[25]	Shen, F	Extension of HPO using embedding of phenotype knowledge resources
<b>Requirement 3: Novel Data for Phenotyping</b>		
[26]	Trace, JM	Using voice to diagnose Parkinson's disease
<b>Requirement 4: Temporal Phenotyping and Subtyping via Similarity Metrics</b>		
[27]	Mate, S	A graphical model of temporal constraints
[28]	Meng, W	Temporal phenotyping of cancer treatment pathways
[29]	Zhao, J	Temporal phenotyping via tensor factorization
[30]	Chen, X	Phenotypic similarity for rare diseases
[31]	Xu, Z	Subtyping for acute kidney injury
<b>Requirement 5: Scalability</b>		
[32]	Zhang, L	Automated grouping of medical codes
[33]	Chen, P	Deep representation learning for phenotyping

challenges in data quality, algorithm portability and interoperability, process efficiency and scalability. Table 1 lists the primary contributions of these twenty articles for deep phenotyping towards the above five topic areas.

## 2. Natural language processing for deep phenotyping

Four articles in this special issue leveraged NLP on different data resources for phenotyping [14–17].

Datta et al. provided a methodology review that provides a frame semantic overview of NLP-based information extraction from EHR notes [14]. Using cancer as an example, this article contributes a model for identifying important disease-specific information using NLP techniques and serves as a useful resource for future researchers requiring disease-specific information extracted from EHR notes.

Liu et al. extracted symptom concepts from clinical notes to stratify patients with mental illness [15]. Contextual terms were extracted to identify constellations of symptoms in a cohort of patients diagnosed with schizophrenia and related disorders. Topic modeling and dimensionality reduction were applied to identify similar groups of patients, who were further evaluated through visualization and interrogation of clinically interpretable weighted features.

Lyudovyyk et al. extracted differentiating clinical phenotypes from pathology reports and combined these phenotypes with biological pathways to identify novel cancer subtypes with prognostic value [16]. This paper shows the value of integrating multi-level biological and clinical observations for deep phenotyping.

Liu et al. addresses the portability challenge facing NLP phenotyping algorithms by presenting an ensemble-based study [17]. Compared to the previously published ensemble methods, this study initially

reported comprehensive comparative effectiveness of four different ensemble techniques over four widely-adopted NLP systems, i.e., MetaMapLite [34], MedLEE [35], ClinPhen [36], and cTAKES [37]. The authors evaluated the performance of different approaches in identifying generic phenotypic concepts and patient-specific phenotypic concepts, respectively, and demonstrated that both tasks benefit from the ensemble techniques.

## 3. Standards: uses and development

A total of eight papers fall into the category of standards use and development.

The first paper leverages the fairly new standard, FHIR (Fast Healthcare Interoperability Resource), for phenotyping. Hong et al. contributed a FHIR-based framework for phenotyping [18] and demonstrated its effectiveness in identifying patients with obesity and multiple comorbidities from clinical text. Given the National Institutes of Health's endorsement of FHIR as the primary data standard for supporting clinical research, this work is timely and relevant.

The next cluster of contributions [19–23] are from the burgeoning open science community, The Observational Health Data Sciences and Informatics (OHDSI) consortium ([www.ohdsi.org](http://www.ohdsi.org)), in collaboration with the eMERGE network. All of these papers use the Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM) to standardizing phenotyping algorithms, data models, or evaluation methods. These standards-based solutions enable transparent collaboration and improve the scalability and efficiency of deep phenotyping.

From firsthand knowledge, retrospective analysis and user surveys, Shang et al. summarized all of the manual effort required to implement electronic phenotypes within the eMERGE Network [19]. This study introduces a novel Knowledge-Interpretation-Programming (KIP) metric to measure portability of phenotype algorithms for quantifying such efforts across the eMERGE Network. The authors concluded that the OMOP CDM can be employed to improve the portability for some 'knowledge-oriented' tasks.

Hripcsak et al. conducted a network-wide study in the eMERGE Network and further demonstrated that the OMOP CDM can facilitate phenotype transfer across network sites and minimize manual implementation [20].

Ostropolets et al. reported lessons learned in adapting EHR-derived phenotypes to claims data, which contain relatively more limited information [21]. The authors succeeded in improving the generalizability and consistency of the chronic kidney disease (CKD) phenotypes by using data and vocabulary standardized by the OMOP CDM. However, performance varied across datasets, implying that even when using a standardized vocabulary, it is important to identify and address coding and data heterogeneity to improve the performance of electronic phenotypes.

Reps et al. developed a novel generalizable probabilistic phenotype model, Current Risk of Smoking Status (CROSS), for current smoking status identification using claims data, which often contains limited data [22]. CROSS can be readily implemented to any US insurance claims mapped to the OMOP CDM and will be useful to impute smoking status when conducting epidemiology studies where smoking is a known confounder but smoking status is not recorded.

A big challenge facing phenotyping research is the lack of rigorous evaluation, which often entails laborious and expensive manual gold standard generation. In order to address this challenge, Swerdel et al. from the OHDSI community contributed an open-source evaluation tool called PheValuator to estimate phenotype algorithm performance [23]. A key contribution of this work is in that it enables scalable, unsupervised evaluation of electronic phenotypes without laborious manual gold standards generation.

Existing standards often have limitations in their content coverage. Warner et al. extended the widely adopted OMOP CDM to expand coverage within the cancer domain by defining a new standard vocabulary for chemotherapy regimen [24]. Similarly, Shen et al. developed a scalable knowledge engineering method to enrich node embedding for HPO [25]. The authors parsed disease-phenotype associations contained in heterogeneous knowledge resources such as OMIM and Orphanet to enrich non-inheritance relationships among phenotypic nodes in HPO.

#### 4. Disease subtyping using novel data sources and temporal reasoning

Trace et al. applied conventional machine learning methods on voice signals to differentiate participants with Parkinson's disease (PD) who exhibit little to no symptoms from healthy controls. They confirmed that voice may serve as a deep phenotype for Parkinson's disease [26].

Both conventional machine learning and emerging deep learning methods have been leveraged for identifying disease subtypes based on temporal characterization of diseases and patient similarity measures.

Mate et al. addressed the temporal complexity of cohort queries and the limitations in existing query tools by creating a novel graphical model for representing temporal cohort queries [27]. This work is a significant applied extension of Allen's time interval algebra. The model was demonstrated to be effective in representing typical temporal phenotype queries using the public MIMIC data.

Meng et al. developed temporal phenotyping methods to derive cancer treatment pathways within a large insurance claims dataset [28]. The authors aggregated lines of therapy information via clustering followed with data visualization to derive temporal cancer phenotypes in support of disease management and progression prediction.

Zhao et al. applied a constrained non-negative tensor-factorization approach on electronic health records to detect temporal phenotypes of complex cardiovascular diseases (CVD) [29]. From a cohort of 12,380 CVD adults, they identified 14 subphenotypes. Through the association analysis with estimated CVD risk for each subtype, they found novel phenotypic topics such as Vitamin D deficiency, depression, and urinary infections. Through a survival analysis, the authors found different risks of subsequent myocardial infarction following the diagnosis of CVD among the six most prevalent topics ( $p < 0.0001$ ), indicating their correspondence to clinically meaningful subphenotypes of CVD. Of note, this study leveraged a coding standard called PheCode [38] to reduce dimensions.

#### 5. Disease subtyping based on patient similarity

Chen et al. developed a comprehensive phenotype similarity metric integrating clinical and questionnaire data for subtyping rare diseases and applied it to ciliopathies [30]. The computed similarity was then validated using genomic data.

Similarly, Lyudovik et al. evaluated the use of genomic test reports ordered for cancer patients in order to derive cancer subtypes and to identify biological pathways predictive of poor survival outcomes [16]. A novel patient similarity metric based on affected biological pathways was proposed. The authors demonstrated that this approach identified subtypes of prognostic value, linked to survival, with implications for precision treatment selection and a better understanding of the underlying disease.

Xu et al. used a memory network-based deep learning approach to discover three acute kidney injury (AKI) sub-phenotypes using EHR data [31]. Group one had an average age of  $63.0 \pm 17.3$  years and mild loss of kidney excretory function, characterizing patients more likely to develop stage I AKI. Group two had an average age of  $66.8 \pm 10.4$  years and severe loss of kidney function, characterizing

patients more likely to develop stage III AKI. Group three had an average age of  $65.1 \pm 11.3$  years and moderate loss of kidney function, characterizing patients more likely to develop stage II AKI.

#### 6. Overcoming challenges in data quality and scalability

Zhang et al. addresses the scalability challenge in developing accurate phenotype algorithms while minimizing manual efforts by contributing a data-driven approach to automate grouping medical terms into clinically relevant concepts by combining multiple up-to-date data sources in an unbiased manner [32]. The proposed method consists of a banding step that leverages the prior knowledge from the existing coding hierarchy, and a combining step that performs spectral clustering on an optimally weighted matrix. The resulting ICD groupings enjoy comparable interpretability and consistency with the current ICD hierarchy.

In contrast to manually creating unbiased estimators of treatment effects, which can be time-consuming and subjective, Chen et al. contributed a scalable method for deep representation learning [33] and applied it for individualized treatment effect estimation using EHR data. The automatically trained representation revealed consistent findings with existing medical knowledge and generated new clinical hypotheses.

#### 7. Summary

We observed the following trends in deep phenotyping research from this set of articles (Table 2). Both the collaborative open science consortia of OHDSI and eMERGE collectively have made significant contributions to deep phenotyping research, with each contributing five or six articles to this special issue, and established best practices for standards-based collaborative phenotyping efforts. The OMOP CDM contributed by OHDSI has demonstrated its promise in facilitating the reusability, efficiency, portability and reproducibility of electronic phenotyping within the eMERGE network. Cancer offers more research challenges and opportunities than other diseases for deep phenotyping. Rare diseases may suffer from limited data but still have a huge need for deep phenotyping. One third of the included papers leveraged the emerging deep learning technologies. Sixty percent of the included papers use formal methods or standards, implying the significant value of data standards in deep phenotyping. About one third of the studies leveraged various NLP methods to include important clinical text in phenotyping. NLP will continue to play an important role in phenotyping research. Six (30%) articles used novel data sources, including clinical text, biological pathways, voice, public knowledge bases, and claims data.

We expect that broad sharing of phenotype definitions and inclusion of diverse data-types in those definitions will continue. We anticipate the next set of challenges to be around including time information in phenotyping as well as in defining the criteria when a phenotype ends. We expect advances in incorporating knowledge—such as physiology and previous evidence—into the phenotype generating process. Another fruitful area of investigation will be automated ways of estimating the portability of phenotype definitions and defining conditions under which porting a definition is unlikely to work. The broad sharing of phenotype definitions that already occur will position us well to pursue these research directions.

#### Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

**Table 2**

Observed trends (NLP: using NLP; STND: using formal standards or semantic knowledge resources such as UMLS; DL: using deep learning technologies; SRC: using novel phenotype data sources; SUB: developing subtypes).

Reference	First author	NLP	STND	DL	SRC	SUB	Disease Focus	eMERGE	OHDSI
[14]	Datta, S	X							
[15]	Liu, Q	X			X				
[16]	Lyudoviyk, O	X			X	X	Cancer	X	
[17]	Liu, C	X		X				X	
[18]	Hong, N		X						
[19]	Shang, N		X					X	X
[20]	Hripcsak, G		X					X	X
[21]	Ostropolets, A		X		X				X
[22]	Reps, J		X		X				X
[23]	Swerdel, J		X						X
[24]	Warner, J		X				Cancer		X
[25]	Shen, F	X	X	X	X				
[26]	Trace, JM				X		Parkinson's		
[27]	Mate, S		X						
[28]	Meng, W						Cancer		
[29]	Zhao, J		X	X			Cardiovascular	X	
[30]	Chen, X		X			X	Rare disease		
[31]	Xu, Z	X		X		X	Acute kidney injury		
[32]	Zhang, L		X	X					
[33]	Chen, P			X					
<b>TOTAL</b>		<b>6</b>	<b>12</b>	<b>6</b>	<b>6</b>	<b>3</b>		<b>5</b>	<b>6</b>

## Acknowledgments

The editors acknowledge funding support from the United States National Institutes of Health grants R01LM009886-10, R01LM006910-19, U01HG008680-05 and R01LM011369-07.

## References

- In Clinical Data as the Basic Staple of Health Learning: Creating and Protecting a Public Good: Workshop Summary. Washington (DC), 2010.
- K.M. Newton, et al., Validation of electronic medical record-based phenotyping algorithms: results and lessons learned from the eMERGE network, *J. Am. Med. Inform. Assoc.* 20 (e1) (2013) e147–e154.
- J. Pathak, A.N. Kho, J.C. Denny, Electronic health records-driven phenotyping: challenges, recent advances, and perspectives, *J. Am. Med. Inform. Assoc.* 20 (e2) (2013) e206–e211.
- P.N. Robinson, Deep phenotyping for precision medicine, *Hum. Mutat.* 33 (5) (2012) 777–780.
- G. Hripcsak, et al., Unlocking clinical data from narrative reports: a study of natural language processing, *Ann. Intern. Med.* 122 (9) (1995) 681–688.
- C.M. Cutillo, et al., Machine intelligence in healthcare—perspectives on trustworthiness, explainability, usability, and transparency, *NPJ Digit Med.* 3 (2020) 47.
- G. Hripcsak, D.J. Albers, High-fidelity phenotyping: richness and freedom from bias, *J. Am. Med. Inform. Assoc.* 25 (3) (2018) 289–294.
- K.B. Waghlikar, et al., Extending i2b2 into a framework for semantic abstraction of EHR to facilitate rapid development and portability of Health IT applications, *AMIA Jt. Summits Transl. Sci. Proc.* 2019 (2019) 370–378.
- C. Sohrabi, et al., World Health Organization declares global emergency: a review of the 2019 novel coronavirus (COVID-19), *Int. J. Surg.* 76 (2020) 71–76.
- H.Y.F. Wong, et al., Frequency and distribution of chest radiographic findings in COVID-19 positive patients, *Radiology* (2019) 201160.
- S.H. Yoon, et al., Chest radiographic and CT findings of the 2019 novel coronavirus disease (COVID-19): analysis of nine patients treated in Korea, *Kor. J. Radiol.* 21 (4) (2020) 494–500.
- J. Lu, et al., Disruptions of the olfactory and default mode networks in Alzheimer's disease, *Brain Behav.* 9 (7) (2019) e01296.
- J. Sun, et al., Supervised patient similarity measure of heterogeneous patient records, *ACM SIGKDD Explorat. Newsletter* 2012 (14) (2012) 16–24.
- S. Datta, E.V. Bernstam, K. Roberts, A frame semantic overview of NLP-based information extraction for cancer-related EHR notes, *J. Biomed. Inform.* 100 (2019) 103301.
- Q. Liu, et al., Symptom-based patient stratification in mental illness using clinical notes, *J. Biomed. Inform.* 98 (2019) 103274.
- O. Lyudoviyk, et al., Pathway analysis of genomic pathology tests for prognostic cancer subtyping, *J. Biomed. Inform.* 98 (2019) 103286.
- C. Liu, et al., Ensembles of natural language processing systems for portable phenotyping solutions, *J. Biomed. Inform.* 100 (2019) 103318.
- N. Hong, et al., Developing a FHIR-based EHR phenotyping framework: a case study for identification of patients with obesity and multiple comorbidities from discharge summaries, *J. Biomed. Inform.* 99 (2019) 103310.
- N. Shang, et al., Making work visible for electronic phenotype implementation: Lessons learned from the eMERGE network, *J. Biomed. Inform.* 99 (2019) 103293.
- G. Hripcsak, et al., Facilitating phenotype transfer using a common data model, *J. Biomed. Inform.* 96 (2019) 103253.
- A. Ostropolets, et al., Adapting electronic health records-derived phenotypes to claims data: lessons learned in using limited clinical data for phenotyping, *J. Biomed. Inform.* (2019) 103363.
- J.M. Reps, P.R. Rijnbeek, P.B. Ryan, Supplementing claims data analysis using self-reported data to develop a probabilistic phenotype model for current smoking status, *J. Biomed. Inform.* 97 (2019) 103264.
- J.N. Swerdel, G. Hripcsak, P.B. Ryan, PheValuator: development and evaluation of a phenotype algorithm evaluator, *J. Biomed. Inform.* 97 (2019) 103258.
- J.L. Warner, et al., HemOnc: a new standard vocabulary for chemotherapy regimen representation in the OMOP common data model, *J. Biomed. Inform.* 96 (2019) 103239.
- F. Shen, et al., HPO2Vec+: leveraging heterogeneous knowledge resources to enrich node embeddings for the Human Phenotype Ontology, *J. Biomed. Inform.* 96 (2019) 103246.
- J.M. Tracy, et al., Investigating voice as a biomarker: deep phenotyping methods for early detection of Parkinson's disease, *J. Biomed. Inform.* (2019) 103362.
- S. Mate, et al., A method for the graphical modeling of relative temporal constraints, *J. Biomed. Inform.* 100 (2019) 103314.
- W. Meng, et al., Temporal phenotyping by mining healthcare data to derive lines of therapy for cancer, *J. Biomed. Inform.* 100 (2019) 103335.
- J. Zhao, et al., Detecting time-evolving phenotypic topics via tensor factorization on electronic health records: cardiovascular disease case study, *J. Biomed. Inform.* 98 (2019) 103270.
- X. Chen, et al., Phenotypic similarity for rare disease: ciliopathy diagnoses and subtyping, *J. Biomed. Inform.* 100 (2019) 103308.
- Z. Xu, et al., Identifying sub-phenotypes of acute kidney injury using structured and unstructured electronic health record data with memory networks, *J. Biomed. Inform.* 102 (2020) 103361.
- L. Zhang, et al., Automated grouping of medical codes via multiview banded spectral clustering, *J. Biomed. Inform.* 100 (2019) 103322.
- P. Chen, et al., Deep representation learning for individualized treatment effect estimation using electronic health records, *J. Biomed. Inform.* 100 (2019) 103303.
- D. Demner-Fushman, W.J. Rogers, A.R. Aronson, MetaMap Lite: an evaluation of a new Java implementation of MetaMap, *J. Am. Med. Inform. Assoc.* 24 (4) (2017) 841–844.
- C. Friedman, et al., A general natural-language text processor for clinical radiology, *J. Am. Med. Inform. Assoc.* 1 (2) (1994) 161–174.
- C.A. Deisseroth, et al., ClinPhen extracts and prioritizes patient phenotypes directly from medical records to expedite genetic disease diagnosis, *Genet. Med.* 21 (7) (2019) 1585–1593.
- G.K. Savova, et al., Mayo clinical Text Analysis and Knowledge Extraction System (cTAKES): architecture, component evaluation and applications, *J. Am. Med. Inform. Assoc.* 17 (5) (2010) 507–513.
- W.Q. Wei, et al., Evaluating phecodes, clinical classification software, and ICD-9-CM codes for phenome-wide association studies in the electronic health record, *PLoS ONE* 12 (7) (2017) e0175508.

Chunhua Weng\*  
Department of Biomedical Informatics, Columbia University, New York, NY,  
USA  
E-mail address: [chunhua@columbia.edu](mailto:chunhua@columbia.edu).  
Nigam H Shah  
Medicine - Biomedical Informatics Research, Stanford University, Stanford,

CA, USA  
E-mail address: [nigam@stanford.edu](mailto:nigam@stanford.edu).  
George Hripcsak  
Department of Biomedical Informatics, Columbia University, New York, NY,  
USA  
E-mail address: [gh13@cumc.columbia.edu](mailto:gh13@cumc.columbia.edu).

---

\* Corresponding author.