



BMJ Open

Ascertaining Framingham heart failure phenotype from inpatient electronic health record data using natural language processing: a multicentre Atherosclerosis Risk in Communities (ARIC) validation study

Carlton R Moore ¹, Saumya Jain ², Stephanie Haas,³ Harish Yadav,³ Eric Whitsel,² Wayne Rosamand,² Gerardo Heiss,² Anna M Kucharska-Newton²

To cite: Moore CR, Jain S, Haas S, *et al.* Ascertaining Framingham heart failure phenotype from inpatient electronic health record data using natural language processing: a multicentre Atherosclerosis Risk in Communities (ARIC) validation study. *BMJ Open* 2021;**11**:e047356. doi:10.1136/bmjopen-2020-047356

► Prepublication history and supplemental material for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2020-047356>).

Received 28 November 2020
Accepted 05 May 2021



© Author(s) (or their employer(s)) 2021. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

Correspondence to

Dr Carlton R Moore;
crmoore@med.unc.edu

ABSTRACT

Objectives Using free-text clinical notes and reports from hospitalised patients, determine the performance of natural language processing (NLP) ascertainment of Framingham heart failure (HF) criteria and phenotype.

Study design A retrospective observational study design of patients hospitalised in 2015 from four hospitals participating in the Atherosclerosis Risk in Communities (ARIC) study was used to determine NLP performance in the ascertainment of Framingham HF criteria and phenotype.

Setting Four ARIC study hospitals, each representing an ARIC study region in the USA.

Participants A stratified random sample of hospitalisations identified using a broad range of International Classification of Disease, ninth revision, diagnostic codes indicative of an HF event and occurring during 2015 was drawn for this study. A randomly selected set of 394 hospitalisations was used as the derivation dataset and 406 hospitalisations was used as the validation dataset.

Intervention Use of NLP on free-text clinical notes and reports to ascertain Framingham HF criteria and phenotype.

Primary and secondary outcome measures NLP performance as measured by sensitivity, specificity, positive-predictive value (PPV) and agreement in ascertainment of Framingham HF criteria and phenotype. Manual medical record review by trained ARIC abstractors was used as the reference standard.

Results Overall, performance of NLP ascertainment of Framingham HF phenotype in the validation dataset was good, with 78.8%, 81.7%, 84.4% and 80.0% for sensitivity, specificity, PPV and agreement, respectively.

Conclusions By decreasing the need for manual chart review, our results on the use of NLP to ascertain Framingham HF phenotype from free-text electronic health record data suggest that validated NLP technology holds the potential for significantly improving the feasibility and efficiency of conducting large-scale epidemiologic surveillance of HF prevalence and incidence.

Strengths and limitations of this study

- The article describes the first study to evaluate performance of natural language processing (NLP) using free-text clinical notes and reports stored in electronic health records to ascertain Framingham heart failure phenotype in multiple regionally dispersed hospitals in the USA with different health systems.
- NLP performances (sensitivity, specificity, positive-predictive value and agreement) are assessed with the reference standard being manual extraction of prespecified information by trained and certified abstractors, using a highly standardised protocol, with quality assurance programmes in place that monitored accuracy, completeness and repeatability of the process.
- The NLP programme used open-source software (clinical Text Analysis Knowledge Extraction System and Python).
- A limitation to the study is that it only includes a subset of hospitalised patients at risk for acute decompensated heart failure based on diagnostic codes (International Classification of Disease, ninth revision) and therefore is not representative of the general hospitalised population.

INTRODUCTION

Since the passage of the Health Information Technology for Economic and Clinical Health Act in 2009,¹ the use of electronic health records (EHRs) in hospital settings has become nearly ubiquitous. Although in 2008, approximately 9% of hospitals were using EHRs, by 2020 the adoption of EHR use among hospitals is approaching 100%.² This creates unprecedented opportunities for researchers to automate the process

of extracting clinical phenotype from patient medical records through electronic search methods.

Scientific organisations and experts promote leveraging electronic data as beneficial to the future of research, public health surveillance and quality improvement initiatives.³ The Working Group on Epidemiology and Population Sciences established by the National Heart, Lung and Blood Institute identified e-epidemiology as a strategic priority for research, with recommendations for studies to ‘determine the validity, reliability and scalability of electronic tools for data collection’.⁴ Clinical phenotypes can be efficiently and accurately extracted from EHRs through the application of algorithms integrating structured data elements such as diagnostic codes, clinical laboratory data and medication lists.⁵ Less well-studied is the use of natural language processing (NLP) of free-text clinical notes stored in EHRs for the ascertainment of complex clinical phenotypes and syndromes.

We focused this study on the use of NLP for the ascertainment of heart failure (HF), a leading cause of hospital admissions and mortality among older adults in the USA.⁶ HF is a complex clinical syndrome characterised by the heart’s inability to supply blood flow sufficient to meet the needs of the body. It is estimated to affect 5.7 million American adults and its prevalence is expected to rise to 8.4 million by 2030.⁷ Reflecting the heterogeneous nature of HF syndromes, there is no universally accepted diagnostic schema for HF that adequately classifies all patients across this syndrome’s pathophysiology, ranging from HF with reduced left ventricular ejection fraction (LVEF) to HF with preserved LVEF (diastolic dysfunction). Signs and symptoms of HF may differ from patient to patient and clinical judgement is typically required to establish a diagnosis of HF for a given patient. The goal of this study is to determine the extent to which accurate EHR-based extraction of Framingham HF criteria phenotypes and HF event classification⁸ can be performed in an automated fashion from clinical notes. We sampled inpatient EHR at four geographically dispersed hospitals with disparate healthcare systems for automated processing and used as a benchmark for our performance an established, standardised protocol of record abstraction and classification.⁷

METHODS

Study population

From 2005 through 2014, the Atherosclerosis Risk in Communities (ARIC) study⁹ conducted community surveillance of HF hospitalisations, classified according to the Framingham schema,⁸ for residents aged 55–84 years in four regions in the USA.^{10 11} To produce annual event rates of HF, eligible hospitalisations from a sample of discharges from acute care hospitals located in ARIC study communities were manually abstracted and events classified according to the presence of the Framingham HF classification criteria.⁸ A hospitalisation was considered eligible for inclusion based on specific primary or

secondary International Classification of Disease, ninth revision, Clinical Modification codes (HF: 428; rheumatic heart disease: 398.91; hypertensive heart disease with congestive heart failure: 402.0, 402.11 or 402.91; hypertensive heart disease and renal failure with HF: 404.01, 404.03, 404.13, 404.91 or 404.93; acute cor-pulmonale: 415.0; chronic pulmonary heart disease, unspecified: 416.9; other primary cardiomyopathies: 425.4; acute oedema of lung, unspecified: 518.4; dyspnoea and respiratory abnormalities: 786.0). Extraction of prespecified information was performed manually by trained and certified abstractors, using a highly standardised protocol, with quality assurance programmes in place that monitored accuracy, completeness and repeatability of the process.¹² A stratified random sample of these hospitalisations occurring during 2015 in four ARIC study hospitals in different ARIC study regions was drawn for the study. A randomly selected set of 394 records was employed as the derivation dataset; the remainder was set aside as the validation set (table 1, N=406). There were no statistically significant differences in patient demographics between the derivation and validation datasets.

Patient and public involvement

HF is a leading cause of hospital admissions and mortality among older adults in the USA.⁶ It is estimated to affect 5.7 million American adults and its prevalence is expected to rise to 8.4 million by 2030.⁷ Therefore, the study outcomes are likely to be a high priority for patients. However, patients were not directly involved in the study design, conduct or outcomes of the research project.

Study design

The primary goal of this study is to determine the accuracy with which EHR-based NLP algorithms can be used to (1) extract Framingham HF criteria variables (table 2) from free-text clinical notes and (2) ascertain the HF phenotype according to the Framingham schema.⁸ As shown in table 2, HF is present if at least two major Framingham criteria are met, or one major and two minor criteria are met. The study also seeks to assess NLP performance reproducibility in ascertainment of Framingham HF phenotype across the four study hospitals.

Data manually extracted by certified ARIC abstractors following a standardised protocol¹³ were used as the reference standard to assess the EHR-based performance of NLP. We used the derivation dataset of 394 records from the four study hospitals to develop the NLP algorithms to extract Framingham HF criteria variables. Once the NLP algorithms were optimised, we assessed NLP performance using a separate validation dataset of 406 unique patient records (table 1).

Figure 1 summarises the study design in which analysis of free-text clinical notes stored in EHRs was compared with manually abstracted Framingham HF phenotype criteria variables (reference standard) from hospitalisations occurring in 2015 at four study hospitals enrolled in the ARIC study (table 1). EHR clinical note types

Table 1 Hospital and patient characteristics for the validation dataset (N=406)

	Hospital ID			
	A	B	C	D
Hospital characteristic				
EHR vendor	Epic	Epic	Epic	Allscripts
Region	South-east	South	North	East
Status	Academic	Academic	Academic	Non-academic
Hospital bed size	873	700	385	247
Abstracted records (N)	122	46	117	121
Patient characteristics				
Age, mean (SD)	73.2 (10.6)	70.4 (10.5)	77.4 (9.9)	78.6 (10.3)
Female, %	43.2	39.1	61.5	59.0
Identified as white, %	54.4	8.7	90.6	93.4
No health insurance, %	0.0	0.0	0.9	9.8
Medicaid insurance, %	4.0	10.9	13.7	1.6

EHR, electronic health record.

used for the analysis included emergency department notes, hospital admission notes, discharge summaries and imaging studies, when available. A structured data element (≥ 4.5 kg weight change during the hospitalisation) was also included.

Extracting HF phenotype criteria from clinical notes in EHRs using NLP

We developed an NLP system using the open-source Apache clinical Text Analysis and Knowledge Extraction Tool¹⁴ (clinical Text Analysis Knowledge Extraction System (cTAKES)) and Python¹⁵ programming software. cTAKES is an NLP programme specifically designed to analyse

free-text clinical notes. It includes specific modules for clinical concept coding and negation status. Concept coding from the Unified Medical Language System¹⁶ was used to identify HF phenotype criteria (table 2), such as ‘paroxysmal nocturnal dyspnoea’, and associate them with standardised concept unique identifiers (CUI), such as ‘C1956415’, that can easily be referenced in a Python programme. The cTAKES programme also assigns a negation status to each concept it identifies in the electronic clinical text. For example, the HF criterion ‘no paroxysmal nocturnal dyspnoea’ is processed by cTAKES by first assigning the CUI ‘C1956415’ to the ‘paroxysmal nocturnal dyspnoea’ concept and then assigning a negation flag to the concept if it identifies predefined negation terms, such as ‘no’ or ‘denies’, associated with the concept. Our study required an additional layer of Python code to identify non-standard documentation of HF criteria (eg, the abbreviation ‘PND’ for ‘paroxysmal nocturnal dyspnoea’), as well as augmented negation so that HF signs and symptoms not described as new or

Table 2 Framingham⁸ HF phenotype criteria variables

Framingham HF phenotype criteria variables	Criteria
4.5 kg weight change over 5 days during hospitalisation	Major
Jugular venous distension	Major
Hepatojugular reflux	Major
Paroxysmal nocturnal dyspnoea	Major
Orthopnea	Major
Pulmonary basilar rales	Major
S3 gallop	Major
Alveolar/pulmonary oedema on chest X-ray	Major
Cardiomegaly on chest X-ray	Major
Lower extremity oedema	Minor
Hepatomegaly	Minor
Dyspnoea—exertion	Minor
Bilateral pleural effusion	Minor

HF is diagnosed if the following are present: (1) two major criteria or (2) one major and two minor criteria.
HF, heart failure.

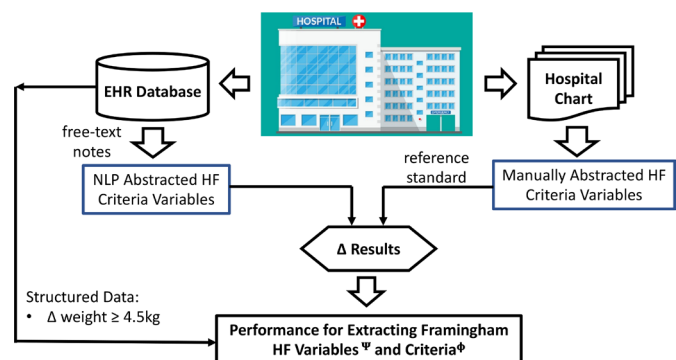


Figure 1 Study design. [‡]Recall (sensitivity), precision (positive-predictive value); [§]recall, precision, Δ estimated HF prevalence, % agreement; EHR, electronic health record; HF, heart failure; NLP, natural language processing.

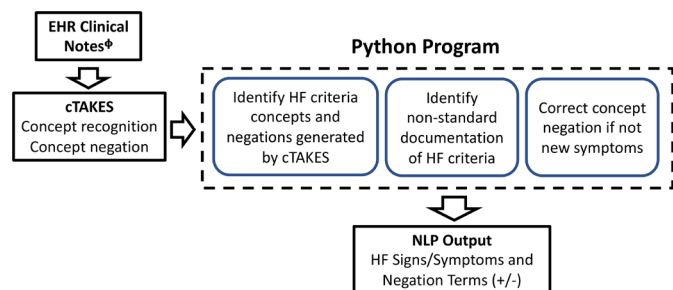


Figure 2 NLP pipeline. [‡]Emergency department notes, hospital admission notes and discharge summaries; cTAKES, clinical Text Analysis Knowledge Extraction System; EHR, electronic health record; HF, heart failure; Negation, determination of whether an HF criteria is negated (eg, patient has oedema vs patient has no oedema); NLP, natural language processing.

worsening were also negated. **Figure 2** shows an overview of the NLP pipeline used to extract Framingham HF phenotype criteria from free-text clinical notes stored in study hospital EHRs. For details of the NLP programme (see online supplemental appendix 1).

Data analysis

We computed sensitivity, specificity and positive-predictive value (PPV) as performance metrics to compare EHR-based HF phenotype criteria with the reference standards (manual review by trained ARIC chart abstractors). Using EHR-based NLP Framingham HF phenotype ascertained criteria (**table 2**), we then calculated the presence or absence of the HF phenotype according to the Framingham⁸ HF schema for the study population, and compared results with Framingham HF phenotype

calculated using manually abstracted Framingham HF criteria from the ARIC study (reference standard). χ^2 and Fisher's exact tests on weighted proportions were used to calculate 95% CIs and p values for EHR-based NLP performance characteristics. All analyses were performed using SAS V.9.4 and Stata/SE V.15.0 software.

RESULTS

EHR performance for extraction of Framingham HF phenotype criteria variables

Table 3 shows the performance of EHR-based NLP abstraction of Framingham HF phenotype criteria from free-text clinical notes, compared with manual chart abstraction for the validation data (see online supplemental appendix 2 for results using derivation data). Cardiomegaly and dyspnoea on exertion showed the best performance at PPV 96.7% and 94.5%, respectively. Conversely, hepatojugular reflux and S3 gallop had the lowest PPVs (0.0% and 11.8%, respectively). A major factor in the poor performance was the low frequency of these variables in the patient sample, 0 and 5 occurrences for hepatojugular reflux and S3 gallop, respectively. Pulmonary oedema demonstrated the best sensitivity (91.7%) and hepatomegaly demonstrated the best specificity (99.0%). See online supplemental appendix 3 for performance of NLP in ascertaining Framingham HF phenotype criteria variables for each study hospital.

NLP performance in the ascertainment of the Framingham HF phenotype from EHR data

Overall, performance of EHR-based ascertainment of Framingham⁸ HF phenotype in the validation dataset was good,

Table 3 NLP performance for abstracting Framingham HF phenotype criteria from EHRs. Validation dataset (N=406)

HF criteria variables (n)*	Sensitivity, % (95% CI)	Specificity, % (95% CI)	PPV, % (95% CI)	Note types used
Weight loss ≥ 4.5 kg† (27)	81.5 (61.9 to 93.7)	96.0 (93.6 to 97.8)	59.5 (43.7 to 75.3)	Structured data
Jugular venous distension (56)	60.7 (46.8 to 73.5)	91.7 (87.6 to 94.8)	61.8 (49.0 to 74.6)	ED, AN
Hepatojugular reflux (0)	N/A	99.7 (98.2 to 100.0)	0.00	ED, AN
PND (27)	55.6 (35.3 to 74.5)	89.4 (85.2 to 92.7)	33.3 (19.2 to 46.7)	ED, AN, DC
Orthopnea (64)	59.4 (46.4 to 71.5)	92.7 (88.7 to 95.6)	67.9 (55.7 to 80.1)	ED, AN, DC
Pulmonary basilar rales (93)	61.3 (50.6 to 71.2)	66.4 (59.7 to 72.6)	43.8 (35.3 to 52.3)	ED, AN, DC
S3 gallop (5)	40.0 (5.3 to 85.3)	95.1 (92.0 to 97.2)	11.8 (0.00 to 27.14)	ED, AN, DC
Pulmonary oedema (48)	91.7 (80.0 to 97.7)	51.0 (44.5 to 57.5)	27.3 (20.4 to 34.2)	ED, AN, DC, IR
Cardiomegaly (162)	54.3 (46.3 to 62.2)	96.0 (90.9 to 98.7)	96.7 (93.0 to 100.0)	ED, AN, DC, IR
Lower extremity oedema (163)	74.8 (67.5 to 81.3)	75.5 (67.7 to 82.2)	77.2 (70.7 to 83.7)	ED, AN, DC
Hepatomegaly (3)	33.3 (0.8 to 90.6)	99.0 (97.2 to 99.8)	33.3 (0.00 to 86.2)	ED, AN, IR
Dyspnoea on exertion (263)	79.1 (73.7 to 83.8)	74.5 (59.7 to 86.1)	94.5 (91.5 to 97.5)	ED, AN, DC
Bilateral pleural effusion (79)	75.9 (65.0 to 84.9)	73.1 (66.5 to 79.0)	51.7 (42.6 to 60.8)	ED, AN, DC, IR

*Instances in total cohort that criteria were identified by manual ARIC abstractors (reference standard).

†Weight loss during hospitalisation based on structured daily patient weight data.

AN, admission note; ARIC, Atherosclerosis Risk in Communities; DC, discharge summary; ED, emergency department; EHRs, electronic health records; HF, heart failure; IR, imaging report; NLP, natural language processing; PND, paroxysmal nocturnal dyspnoea; PPV, positive-predictive value.

Table 4 Performance of NLP-based ascertainment of Framingham HF phenotype

Sensitivity, % (95% CI)	Specificity, % (95% CI)	PPV, % (95% CI)	Agreement, % (95% CI)
78.8 (72.8 to 83.9)	81.7 (75.2 to 87.0)	84.4 (79.5 to 89.3)	80.0 (75.8 to 83.8)

Note types: emergency department notes, hospital admission notes and discharge summaries. HF, heart failure; NLP, natural language processing; PPV, positive-predictive value.

with 78.8%, 81.7%, 84.4% and 80.0% as sensitivity, specificity, PPV and agreement metrics, respectively (table 4).

Performance of NLP-based ascertainment of Framingham HF phenotype across study hospitals

Figure 3 shows EHR-based performance in the ascertainment of Framingham HF phenotype for each of the four study hospitals. There was good reproducibility of NLP performance and no meaningful differences in NLP performance across hospitals for the three performance measures of sensitivity, specificity and agreement (all 95% CIs overlap between hospitals for each performance measure).

DISCUSSION

Here, we report on the derivation and validation of an open-source software NLP application that uses EHR data to ascertain HF according to the established Framingham schema in patients hospitalised in dispersed regions of the USA. EHR-based identification of the Framingham HF phenotype had very good performance characteristics (sensitivity: 78.8%, specificity: 81.7%, PPV: 84.4% and agreement: 80.0%) and was reproducible across the four study hospitals.

Several studies have investigated the use of billing codes and lab results to ascertain Framingham HF phenotype in inpatient settings within single healthcare systems.^{17 18} To our knowledge, this is the first study to describe the performance of EHR-based NLP tools to ascertain Framingham HF phenotype in inpatients from

multiple geographically diverse hospitals from different healthcare systems. Our results compare favourably with studies using ICD-9, diagnosis related group (DRG) codes and lab results to ascertain Framingham HF phenotype. Using ICD-9 and DRG codes, Presley *et al*¹⁷ ascertained Framingham HF phenotype for hospitalised patients in the Veterans Administration (VA) healthcare system.¹⁷ The VA study demonstrated sensitivity of 45.1%, specificity of 99.4% and a PPV of 89.7% for Framingham HF phenotype in population that was homogenous with respect to gender (98.8% male). Using ICD-9 codes, HF medications and lab results, Tison *et al*¹⁸ ascertained Framingham HF phenotype for inpatients within a single healthcare system in Minnesota. Of the multiple study algorithms used in the study, the one with the highest PPV (86.5%) had a sensitivity of only 41.6%.

Our study adds to the growing body of evidence which suggests that NLP has the potential to improve the cost-effectiveness and timeliness of phenotyping in clinical and epidemiological studies by reducing the need for manual chart abstraction.

In this first step towards the development of a robust protocol for EHR-based NLP surveillance of hospitalised HF patients, we designed a prototype system that had good performance in ascertaining Framingham HF phenotype that was reproducible across four hospitals selected to be geographically dispersed. Underlying this reproducibility, however, was considerable effort required to harmonise a single NLP algorithm that accurately and consistently performed well (figure 3) across the four hospitals.

Evaluation of our results revealed several lessons learnt in the extraction of HF phenotype criteria. First, having complete sets of clinical note and report types from hospitals likely had a significant impact on performance. Our study used NLP to process emergency department notes, admission notes, discharge summaries and imaging study reports. Given the notable lack of standardisation of note type nomenclature across hospitals, we found significant variability between the four study hospitals in nomenclature used to identify specific clinical note types. For example, participating hospitals designated discharge summaries as 'Discharge Summaries', 'Discharge PN', 'PMNDIS' and instances in which the discharging physicians name was concatenated with 'Discharge PN' (eg, Smith Discharge PN). To properly capture phenotypes and clinical outcomes from EHRs requires overcoming a lack of standardised nomenclature, variability in standards for defining and recording data elements, and uncertain collection of longitudinal information or data

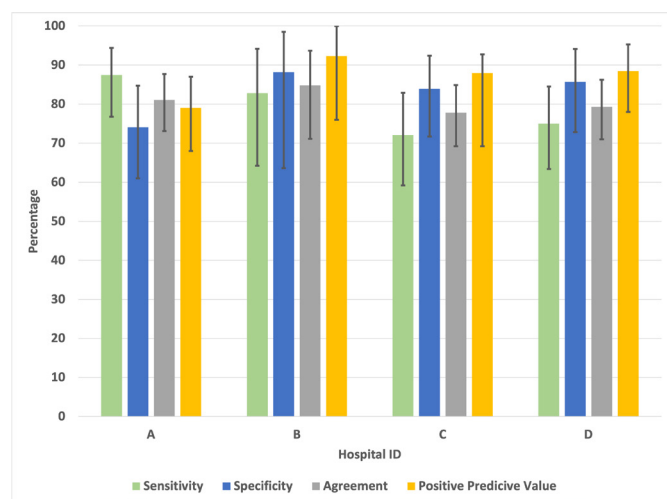


Figure 3 EHR-based performance for Framingham HF phenotype by hospital with 95% CIs. EHR, electronic health record; HF, heart failure.

across settings of care. In contrast, these are all features embedded in the standardised community surveillance registry that systematically gathers data entered by many clinicians in numerous hospitals, and served as the benchmark to validate our HF phenotype identification and event classification from EHR. As is typically the case for dedicated registries, ARIC's data element extraction from records is performed by trained abstractors according to specific definitions, standardised procedures, and use of specialised forms leading to highly reliable and valid information under quality control monitoring. Data in EHRs by contrast are captured in the process of patient care by various members of the clinical team, for purposes other than event ascertainment or analysis. Although several efforts exist to establish common data models for EHR data,^{19 20} such models are not yet in widespread use and standardised definitions when documenting patient care are uncommon.

The second lesson learnt from the study was the challenge in optimising NLP performance to accurately determine negation for Framingham HF phenotype criteria variables documented in clinical notes. We observed multiple instances in which clinicians documented negative HF signs and symptoms phrased as 'patient denies cough, fever, abdominal pain, chest pain, dyspnoea'. In this example, it was often difficult to accurately assess whether an HF phenotype criteria variable was negated by 'denies'. Similarly, formatting of negation terms often varied by clinician and hospital and included terms such as 'no', 'denies', 'negative', 'neg', '(-)', '-', 'patient does not report'; among other idiosyncratic terminology. Another challenge was establishing negation when clinicians described conditions in discharge summaries under which it was appropriate for patients to take a given medication. For example, 'use albuterol inhaler four times daily as needed for dyspnoea'. In this case, the 'dyspnoea' Framingham HF criterion should be negated because the patient is not currently experiencing dyspnoea, a conditional symptom in which a particular medication should be used.

There are limitations to our study results. The study population represents a sample of hospitalised patients selected for the likelihood of having congestive HF based on ICD-9-CM codes (the prevalence of Framingham HF was 52.0% for NLP and 55.8% for manual chart abstraction). However, this limitation can be mitigated by automated screening of patients using the same ICD-9-CM codes before using NLP ascertainment of Framingham HF criteria. Generalisability of study findings to other populations has not been tested. Furthermore, among the metrics used to ascertain NLP performance, estimated PPV is influenced by the prevalence of the condition. Lastly and not unexpectedly, PPV performed poorly for Framingham HF criteria that occurred infrequently in the patient population. Examples of those were hepatjugular reflux (n=0/406), hepatomegaly (n=3/406), S3 gallop (n=5/406) and PND (n=27/406) had PPVs of 0.0%, 33.3%, 11.8% and 33.3%; respectively (table 3).

Nonetheless, because of their low prevalence in the study population, these criteria likely had a relatively small impact on the determination of the Framingham HF phenotype prevalence.

The means to assess the population burden of HF and the impact of medical interventions and public health policies on these metrics are limited, and largely rely on efforts by professional organisations such as the American Heart Association²¹ drawn from various NIH-supported observational studies. Our data suggest that NLP has good performance characteristics in determining Framingham HF phenotype in hospitals from four distinct regions of the country. Such estimates do not substitute for comprehensive population data, nor are they regionally (or nationally) representative, and they do not lend themselves to estimation of population burden metrics or temporal trends. A 2011 report from the Institute of Medicine²² recommended a national surveillance programme to be put in place funded by the Affordable Care Act,²² but questions persist about the feasibility of community surveillance that can efficiently incorporate EHR capabilities for accurate estimates of disease burden and to monitor trends in cardiovascular diseases. To accomplish this, such surveillance should be able to link EHR resources to population denominators, harmonise diverse EHRs and implement information extraction tools of known validity and portability, while safeguarding patient privacy and be robust to changes in diagnostic fashion, technologies and coding practices. Such challenges need careful attention to realise the potential of EHR-enabled community surveillance. The alternative—the current inability to monitor population burden and trends—represents a significant impediment to the ability to gauge the impact of health care and public health initiatives on the burden of, and trends in the most prominent contributors to morbidity, mortality and healthcare expenditures in the USA.

Importantly, the lack of community surveillance programmes encumbers the progress in the understanding of and in reducing health disparities in the incidence of the major cardiovascular health events and their outcomes.²³ Regional epidemiologic surveillance programmes, such as ARIC's, indicate that during the years 2005–2012, annual rates of incident hospitalised HF increased in all race–gender groups, but markedly so for black women. Ongoing HF surveillance efforts are therefore needed to identify vulnerable population subgroups and develop effective prevention strategies.

Future directions for our project include developing a user-friendly interface to adjust NLP algorithms based on institution-specific patterns in documentation of negations, as well as investigating the use of machine-learning technology to optimise performance of the current rule-based NLP system.²⁴ Specifically, our goal is to approach 100% sensitivity while optimising specificity and PPV.

In conclusion, by decreasing the need for manual chart review, our results on the use of NLP to ascertain Framingham HF phenotype from free-text EHR data suggest

that validated NLP technology holds the potential for significantly improving the feasibility and efficiency of conducting large-scale epidemiologic surveillance of HF prevalence and incidence.

Author affiliations

¹Medicine, University of North Carolina at Chapel Hill School of Medicine, Chapel Hill, North Carolina, USA

²Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

³School of Information and Library Science, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA

Acknowledgements The authors thank the staff and participants of the ARIC study for their important contributions.

Contributors All named authors are responsible for the reported research. CRM, EW, WR, SH, GH and AMK-N: substantially contributed to the conception of the project, acquisition of data, interpretation of data for the study and writing of the manuscript. CRM, SJ and HY: contributed significantly to data analysis for the study.

Funding The Atherosclerosis Risk in Communities study has been funded in whole or in part with federal funds from the National Heart, Lung and Blood Institute, National Institutes of Health, Department of Health and Human Services, under contract nos (HHSN2682017000011, HHSN2682017000021, HHSN2682017000031, HHSN2682017000051 and HHSN2682017000041).

Competing interests None declared.

Patient consent for publication Not required.

Ethics approval This study was approved by the University of North Carolina Institutional Review Board (IRB Study ID#: 19-3462). All participants gave informed consent before taking part in the study.

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement Data are available upon reasonable request. All raw data for the study will be available upon request.

Supplemental material This content has been supplied by the author(s). It has not been vetted by BMJ Publishing Group Limited (BMJ) and may not have been peer-reviewed. Any opinions or recommendations discussed are solely those of the author(s) and are not endorsed by BMJ. BMJ disclaims all liability and responsibility arising from any reliance placed on the content. Where the content includes any translated material, BMJ does not warrant the accuracy and reliability of the translations (including but not limited to local regulations, clinical guidelines, terminology, drug names and drug dosages), and is not responsible for any error and/or omissions arising from translation and adaptation or otherwise.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

ORCID iDs

Carlton R Moore <http://orcid.org/0000-0001-7732-4829>

Saumya Jain <http://orcid.org/0000-0002-4083-3918>

REFERENCES

- Bowes WA. Assessing readiness for meeting meaningful use: identifying electronic health record functionality and measuring levels of adoption. *AMIA Annual Symposium proceedings / AMIA Symposium AMIA Symposium*, 2010:66–70.
- et alHenry J, Pylypchuk Y, Searcy T. Adoption of electronic health record systems among U.S. Non-Federal acute care hospitals, 2016. Available: <https://dashboard.healthit.gov/evaluations/data-briefs/non-federal-acute-care-hospital-ehr-adoption-2008-2015.php> [Accessed 28 Jan 2020].
- IOM. *To err is human: building a safer health system*. Washington, DC: Institute of Medicine, 2000.
- Boerwinkle E, Crapo JD, Douglas PS. *Strategic transformation of population studies*. NHLBI Advisory Council, 2014.
- Ritchie MD, Denny JC, Crawford DC, et al. Robust replication of genotype-phenotype associations across multiple diseases in an electronic medical record. *Am J Hum Genet* 2010;86:560–72.
- Virani SS, Alonso A, Benjamin EJ, et al. Heart disease and stroke Statistics-2020 update: a report from the American heart association. *Circulation* 2020;141:e139–596.
- Heidenreich PA, Albert NM, Allen LA, et al. Forecasting the impact of heart failure in the United States: a policy statement from the American heart association. *Circ Heart Fail* 2013;6:606–19.
- Ho KK, Pinsky JL, Kannel WB, et al. The epidemiology of heart failure: the Framingham study. *J Am Coll Cardiol* 1993;22:6A–13.
- The Atherosclerosis risk in communities (ARIC) study: design and objectives. The ARIC Investigators. *Am J Epidemiol* 1989;129:687–702.
- White AD, Folsom AR, Chambless LE, et al. Community surveillance of coronary heart disease in the Atherosclerosis risk in communities (ARIC) study: methods and initial two years' experience. *J Clin Epidemiol* 1996;49:223–33.
- Rosamond WD, Chang PP, Baggett C, et al. Classification of heart failure in the Atherosclerosis risk in communities (ARIC) study: a comparison of diagnostic criteria. *Circ Heart Fail* 2012;5:152–9.
- Loehr LR, Agarwal SK, Baggett C, et al. Classification of acute decompensated heart failure: an automated algorithm compared with a physician reviewer panel: the Atherosclerosis risk in Communities study. *Circ Heart Fail* 2013;6:719–26.
- ARIC. Surveillance of heart failure manual of operation. ARIC coordinating center, 2011. Available: https://sites.csc.unc.edu/aric/sites/default/files/public/manuals/Man3a_updated_6172011_0.pdf [Accessed 26 Aug 2020].
- Savova GK, Masanz JJ, Ogren PV, et al. Mayo clinical text analysis and knowledge extraction system (cTAKES): architecture, component evaluation and applications. *J Am Med Inform Assoc* 2010;17:507–13.
- Python. Python. python software Foundation, 2016. Available: <https://www.python.org/doc/> [Accessed 29 Jan 2019].
- NIH/NLM. Unified medical language system (UMLS). National library of medicine, 2016. Available: <https://www.nlm.nih.gov/research/umls/> [Accessed 2 Jun 2016].
- Presley CA, Min JY, Chipman J, et al. Validation of an algorithm to identify heart failure hospitalisations in patients with diabetes within the Veterans health administration. *BMJ Open* 2018;8:e020455.
- Tison GH, Chamberlain AM, Pletcher MJ, et al. Identifying heart failure using EMR-based algorithms. *Int J Med Inform* 2018;120:1–7.
- Stang PE, Ryan PB, Racoosin JA, et al. Advancing the science for active surveillance: rationale and design for the observational medical outcomes partnership. *Ann Intern Med* 2010;153:600–6.
- ONC, Interoperability. Office of the National coordinator for health information technology (onc), 2019. Available: <https://www.healthit.gov/topic/interoperability> [Accessed 17 Sep 2020].
- Writing Group Members, Mozaffarian D, Benjamin EJ, et al. Heart disease and stroke Statistics-2016 update: a report from the American heart association. *Circulation* 2016;133:e38–60.
- IOM. *The future of nursing: leading change, advancing health*. Washington, DC: National Academy of Sciences, 2011.
- Nih morbidity and mortality chart book, 2012. Available: <http://www.nhlbi.nih.gov/research/reports/2012-mortality-chart-book>
- Pradhan S, Elhadad N, South BR, et al. Evaluating the state of the art in disorder recognition and normalization of the clinical narrative. *J Am Med Inform Assoc* 2015;22:143–54.

1 Bowes WA. Assessing readiness for meeting meaningful use: identifying electronic health record functionality and measuring