

# BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email [info.bmjopen@bmj.com](mailto:info.bmjopen@bmj.com)

# BMJ Open

## Real-World Evidence in Heart Failure: Data Quality and Evidence Validity

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2023-073178
Article Type:	Original research
Date Submitted by the Author:	01-Mar-2023
Complete List of Authors:	Garan, Arthur Reshad; Harvard Medical School Monda , Keri; Amgen Inc Dent-Acosta, Ricardo ; Amgen Inc Riskin , Daniel; Verantos Gluckman, Ty; Providence St Joseph Health
Keywords:	Heart failure < CARDIOLOGY, Cardiac Epidemiology < CARDIOLOGY, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

## Real-World Evidence in Heart Failure: Data Quality and Evidence Validity

**Short title:** Real-World Evidence in Heart Failure

**Authors:** A. Reshad Garan, MD, MS; Keri L. Monda, Ph.D.; Ricardo E. Dent-Acosta; MD; Dan Riskin, MD; Ty J. Gluckman, MD, MHA

**Affiliations:** Department of Medicine, Division of Cardiology, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, MA, USA (A.R.G.). The Center for Observational Research and Medical Affairs, Amgen Inc., Thousand Oaks, CA, USA (K.L.M., R.E.D.-A.). Department of Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA (K.L.M.). Verantoss, Inc., Menlo Park, CA, USA (D.R.). Department of Surgery, Stanford University School of Medicine, Stanford, CA, USA (D.R.). Center for Cardiovascular Analytics, Research and Data Science (CARDS), Providence Heart Institute, Providence St. Joseph Health, Portland, OR, USA (T.J.G.)

**Address for correspondence:** A. Reshad Garan, MD, MS, Department of Medicine, Division of Cardiology, Beth Israel Deaconess Medical Center, 185 Pilgrim Road Deaconess 3, Boston, MA 02215, USA. Phone: 617/632-7737. Fax: 617/632-7620.

Email: [agaran@bidmc.harvard.edu](mailto:agaran@bidmc.harvard.edu)

**Word count:** 2673 words

**Keywords:** artificial intelligence; heart failure; phenotype; real-world evidence; electronic health record.

### Abstract

1  
2  
3 **Objective:** Quantitatively evaluate the quality of data underlying real-world evidence  
4 (RWE) in heart failure (HF).  
5  
6

7 **Design:** Retrospective comparison of accuracy in identifying HF patients and  
8 phenotypic information was made using traditional (i.e., structured query language  
9 applied to structured EHR data) and advanced (i.e., AI applied to unstructured EHR  
10 data) RWE approaches. The performance of each approach was measured by the  
11 harmonic mean of precision and recall (F1 score) using manual annotation of medical  
12 records as a reference standard.  
13  
14  
15  
16  
17  
18  
19  
20

21 **Setting:** EHR data from a large academic healthcare system in North America between  
22 2015 and 2019, with an expected catchment of approximately 500,000 patients.  
23  
24  
25

26 **Population:** 4288 encounters for 1155 patients aged 18 to 85 years, with 472 patients  
27 identified as having HF.  
28  
29

30 **Outcome measures:** HF and associated concepts, such as comorbidities, left  
31 ventricular ejection fraction, and selected medications.  
32  
33  
34  
35

36 **Results:** The average F1 scores across 19 HF-specific concepts were 49.0% and  
37 94.1% for the traditional and advanced approaches, respectively ( $P < 0.001$  for all  
38 concepts with available data). The absolute difference in F<sub>1</sub> score between approaches  
39 was 45.1% (98.1% relative increase in F<sub>1</sub> score using the advanced approach). The  
40 advanced approach achieved superior F1 scores for HF presence, phenotype, and  
41 associated comorbidities. Some phenotypes, such as HFpEF, revealed dramatic  
42 differences in extraction accuracy based on technology applied, with a 4.9% F<sub>1</sub> score  
43 when using natural language processing (NLP) alone and a 91.0% F<sub>1</sub> score when using  
44 NLP plus AI-based inference.  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

1  
2  
3 **Conclusions:** A traditional RWE generation approach resulted in low data quality in HF  
4 patients. While an advanced approach demonstrated high accuracy, the results varied  
5 dramatically based on extraction techniques. For future studies, advanced approaches  
6 and accuracy measurement may be required to ensure data are fit-for-purpose.  
7  
8  
9  
10  
11  
12

### 13 **Strengths and limitations of this study**

- 14 • Using RWE for HF patients requires demonstrating that the data source and  
15 technologies result in accurate data.
  - 16 • Natural language processing alone lacked context from the longitudinal record,  
17 limiting phenotype identification and study validity.
  - 18 • Findings suggest that advanced methods can enable high-validity RWE for heart  
19 failure patients.
  - 20 • The use of data from a single healthcare system may limit generalizability to  
21 other populations.
- 22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## INTRODUCTION

Heart failure (HF) is a major public health problem with significant associated morbidity, mortality, and cost.<sup>1,2</sup> Despite the availability of novel drugs and devices, morbidity and mortality in HF rivals many malignancies, with a 5-year survival rate as low as 50%.<sup>3-8</sup>

Randomized controlled trials (RCTs) have traditionally been used to assess the safety and efficacy of new therapies and represent a cornerstone for regulatory approval.

However, RCTs are frequently conducted in highly selected populations, typically younger, healthier, and less diverse than patients treated in clinical practice.

Furthermore, such trials often include patients with an established HF diagnosis, receiving guideline-directed medical therapy at tertiary centers, and may not represent the broader HF population. In contrast, registry data usually offers additional insights into more inclusive populations. Even with this, there is potential bias based on inclusion and exclusion criteria. Because HF is a clinically heterogeneous syndrome with numerous etiologies and phenotypes, studying this population can be particularly difficult.

Real-world evidence (RWE) has held promise as a potential means to assess therapeutic benefit outside of clinical trials, with sufficient power to characterize therapeutic impact in HF subgroups. Accordingly, RWE can complement RCTs, extending the findings to patient populations that may have been excluded from or insufficiently enrolled in pivotal trials. To accelerate these and similar precision medicine goals, the 21st Century Cures Act was passed in 2016, which required the United States Food and Drug Administration (FDA) to develop guidance supporting the use of

1  
2  
3 RWE in new drug indications and post-marketing surveillance.<sup>9</sup> In addition, payors have  
4 increasingly utilized RWE to inform reimbursement decisions and are increasingly  
5 demanding credible evidence.<sup>10</sup>  
6  
7  
8  
9

10  
11 Not surprisingly, the quality of RWE hinges on how well real-world data are collected,  
12 processed<sup>11</sup>, and used to inform study questions. Such is the case in HF, where  
13 accurate identification of patients in administrative and other structured data sets is an  
14 ongoing focus.<sup>12-14</sup> Artificial intelligence (AI) applied to unstructured data represents a  
15 novel method of analyzing the electronic health record (EHR). Because of the  
16 importance of data reliability in RWE and the potential to use unstructured data to  
17 achieve data enrichment<sup>15</sup>, we sought to better understand differences in accuracy  
18 between traditional RWE methods and advanced AI approaches for a range of HF-  
19 specific data elements.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

## 32 **METHODS**

33  
34  
35  
36 Varied data sources and applied technologies were used to assess data reliability in  
37 patients with risk factors for HF. Leveraging manual chart abstraction as the reference  
38 standard, comparisons were made between the two methods. The first method used  
39 structured EHR data (e.g., diagnosis codes and problem lists) and standard query  
40 techniques, defined as the 'traditional approach'. The second used unstructured EHR  
41 data (e.g., narratives from primary care and specialty notes) and AI techniques,  
42 described as the 'advanced approach' (Figure 1). The primary objective was  
43 measurement of the accuracy of identified HF-specific elements using traditional and  
44 advanced approaches. We hypothesized that the advanced approach would better  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57



1  
2  
3 identify key HF-specific elements than the traditional approach. Data were deidentified  
4  
5 before study initiation, and the study was determined not to be human subjects  
6  
7 research. Both natural language processing (NLP) and machine-learned inference  
8  
9 technologies used in the advanced approach were provided by Verantos, Inc. (Menlo  
10  
11 Park, CA, USA). The core of AI is a deterministic NLP layer. This layer is built on top of  
12  
13 the GATE NLP architecture.<sup>16</sup> The architecture is used to construct a flexible pipeline for  
14  
15 processing incoming text against English language syntactical rules augmented with a  
16  
17 lexicon based on a clinical vocabulary. The AI-based inference was applied during data  
18  
19 processing. Millions of machine-learned and manually curated associations enable  
20  
21 disambiguation and identification of clinically relevant concepts. As an example of AI-  
22  
23 based inference, a patient with HF on the problem list and a narrative encounter  
24  
25 describing “EF 60%” would not be interpreted by NLP as having HF with preserved  
26  
27 ejection fraction (HFpEF) since the text does not have sufficient information to identify  
28  
29 this condition. On the other hand, AI-based inference would infer HFpEF based on  
30  
31 disparate information in the record.  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## EHR Data Source and Processing

EHR data from primary care encounters between 2011 and 2018 were deidentified and securely transferred to a cloud-based server for analysis. The data set consisted of both structured data (e.g., medical conditions, procedures performed, medications, and problem lists) and unstructured data (e.g., narrative notes from primary care providers and specialists, telephone visits, and other narrative text) (Figure 1).

As the study aimed to test the accuracy of different RWE approaches and not treatment effectiveness, the cohort was enriched for patients with suspected HF based on comorbidities and medications. Specifically, the following filters were applied: records containing both narrative and structured components; narrative length 1,000 characters or more; and at least one of the following problems or medications in structured or unstructured data: myocardial infarction, congestive heart failure, or carvedilol.

A prespecified set of clinical concepts pertinent to patients with HF was extracted using traditional and advanced techniques (Table 1). Problem lists were mapped to Systematized Nomenclature of Medicine (SNOMED) ontology, and unadjudicated claims were mapped to ICD-10 codes. Standard sets of individual codes were used to represent each concept. With the advanced approach, inference incorporating pattern recognition was utilized to identify potentially missing or ignored concepts within the text (e.g., HF being likely in patients with dyspnea and pitting edema on a diuretic).

Specifically, no narrative coding took place before the AI algorithm was used; instead, it was applied directly to the narrative text and then mapped by the algorithm to the SNOMED ontology. Next, manual chart abstraction using the same SNOMED code set

1  
2  
3 was used as a reference to assess the accuracy of the coding by the AI algorithm.  
4  
5 Engineers were blinded to validation data and its corresponding chart abstraction.  
6  
7

## 8 9 **Study End Points and Statistical Analysis**

10  
11  
12 The primary endpoint was the  $F_1$  score for traditional and advanced approaches. The  $F_1$   
13  
14 score is an accuracy measure that combines recall and precision; more specifically, it is  
15  
16 the weighted harmonic mean of these two measures. Secondary endpoints were recall  
17  
18 (i.e., the proportion of patients correctly identified as having the condition, akin to  
19  
20 sensitivity) and precision (i.e., the proportion of patients with HF and its subtypes  
21  
22 correctly identified divided by the total number of patients identified in each cohort akin  
23  
24 to positive predictive value)<sup>17,18</sup> for the traditional and advanced approaches. The  
25  
26 reference standard used to evaluate accuracy of the traditional and advanced  
27  
28 approaches was manual chart abstraction. For each encounter, two independent clinical  
29  
30 annotators labeled each concept and all metadata for that concept. Annotators were  
31  
32 blinded to each other's annotations, and inter-rater agreement was measured by  
33  
34 Cohen's kappa score. Further description of the reference standard methodology is  
35  
36 provided in the Supplemental Material. Results were summarized using descriptive  
37  
38 statistics, and percentages were calculated for categorical variables. Differences in  $F_1$   
39  
40 scores between traditional and advanced approaches were analyzed using the chi-  
41  
42 square test; associated  $P$ -values were reported.  
43  
44  
45  
46  
47  
48  
49

## 50 **RESULTS**

1  
2  
3 A total of 4288 encounters for 1155 patients were examined, of which 472 patients with  
4 HF were identified. Of these, 382 had HF with reduced ejection fraction (HFrEF), 35 had  
5 HF with mildly reduced ejection fraction (HFmrEF), and 55 had HF with preserved  
6 ejection fraction (HFpEF). The reference standard Cohen's kappa score was 0.95,  
7 suggesting high validity.  
8  
9

10  
11  
12  
13  
14  
15  
16 Supplementary Table 1 reports the  $F_1$  score, recall, and precision results achieved with  
17 both approaches. Figure 2 graphically presents  $F_1$  scores for HF diagnoses and Figure  
18 3 includes  $F_1$  scores for symptoms, medications, and comorbid conditions. Overall,  
19 accuracy was significantly greater for the advanced approach (AI applied to  
20 unstructured EHR data) than for the traditional approach (structured query language  
21 applied to structured EHR data) (Supplementary Table 1; Figure 2; Figure 3), with an  
22 absolute difference of 45.1%.  
23  
24  
25  
26  
27  
28  
29  
30

31  
32  
33 With the traditional approach, recall for any HF diagnosis was 46.9% (i.e., 53.1% of  
34 patients with HF were missed entirely) and precision was 95.4%, resulting in an  $F_1$   
35 score of 62.9% ( $P < 0.001$ ). In contrast, with the advanced approach, recall for any HF  
36 diagnosis was 96.0% and precision was 94.7%, resulting in an  $F_1$ -score of 95.3%  
37 ( $P < 0.001$  when  $F_1$  scores for the two approaches were compared) (Supplementary  
38 Table 1; Figure 2). Among HF phenotypes, recall with the advanced approach was  
39 highest with HFrEF, followed by HFpEF and HFmrEF; precision was 100% for all  
40 phenotypes. With the traditional approach,  $F_1$  scores could not be calculated for HFrEF,  
41 HFmrEF, and HFpEF because only less granular HF codes were used (Supplementary  
42 Table 1).  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Accuracy in identifying left ventricular ejection fraction (LVEF) was similarly high with  
4 the advanced approach, with an  $F_1$  score of 96.7%. Data could not be extracted for  
5 LVEF with the traditional approach because no such codes were available within the  
6 EHR, nor did a mechanism to encode LVEF within the problem list or unadjudicated  
7 claims exist (Supplementary Table 1; Figure 2).  
8  
9  
10  
11  
12  
13

14  
15  
16 Accurate identification of HF symptoms was greater with the advanced approach  
17 ( $P<0.001$ ) (Supplementary Table 1; Figure 3A). Whereas identification of commonly  
18 prescribed HF medications was high with both approaches (Supplementary Table 1;  
19 Figure 3B), identification of cardiovascular comorbidities was higher in all cases with the  
20 advanced approach ( $P<0.001$ ) (Supplementary Table 1; Figure 3C).  
21  
22  
23  
24  
25  
26  
27

28 Data concept extraction with the advanced approach greatly depended upon the  
29 technology used. For example, NLP, which ends at the sentence boundary, was only  
30 able to identify HFpEF with an  $F_1$  score of 4.9% because "HFpEF" or "heart failure with  
31 preserved ejection fraction" was rarely written. Conversely, inference, which can find  
32 related items from the longitudinal record, was able to identify both "HF" and "normal  
33 ejection fraction" as separate annotations for HFpEF with an  $F_1$ -score of 91.0%  
34 (Supplementary Table 1; Figure 2).  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

## 45 **DISCUSSION**

46  
47  
48  
49 The utilization of RWE has grown substantially in recent years, driven in part by its  
50 perceived value by clinicians, regulators, and payors. As RWE is increasingly used to  
51 refine care standards through clinical, regulatory, and reimbursement pathways, its  
52  
53  
54  
55  
56  
57

1  
2  
3 accuracy has come under increased scrutiny. This is particularly important for complex  
4 medical conditions, such as HF. Accordingly, we used chart abstraction to quantitatively  
5 evaluate traditional and advanced approaches to define HF-specific data elements. This  
6 allowed us to rigorously evaluate whether commonly used techniques are sufficiently  
7 accurate for observational studies, comparative effectiveness research, and post-  
8 approval safety studies.  
9

10  
11 In this study, we demonstrated that: 1) the use of an advanced, AI-based approach  
12 consistently identified HF phenotypes (i.e., HF<sub>r</sub>EF, HF<sub>m</sub>rEF, and HF<sub>p</sub>EF) more  
13 accurately than a traditional approach; 2) common HF symptoms and comorbid  
14 conditions were consistently and accurately identified using an advanced approach; and  
15 3) medications for HF were accurately identified using both advanced and traditional  
16 approaches. While studies have previously leveraged an AI-based approach to identify  
17 patients with HF,<sup>19-22</sup> our study highlights the discrepancy between traditional EHR  
18 query methods and an AI-based approach standardized against a manual reference.  
19 Given that the accuracy of the data set and appropriateness of the applied technology  
20 are not tested in many RWE studies, there is a high potential for error.<sup>23,24</sup> The current  
21 findings highlight this while also reinforcing the impact that specific AI technologies  
22 (e.g., NLP vs. NLP plus inference) can have on phenotype generation and study  
23 validity.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

49 Accurate phenotyping is paramount in any RWE study that includes HF patients. With  
50 varying etiologies and multiple phenotypes, HF is a clinically diverse syndrome, with  
51 outcomes that may vary between subgroups. In addition, HF patients may have different  
52  
53  
54  
55  
56  
57

1  
2  
3 trajectories, highlighting some of the limitations of using structured data. For example,  
4  
5 LVEF may fluctuate throughout a patient's disease course, with some patients  
6  
7 experiencing recovery of their LVEF with the use of guideline-directed medical therapy.  
8  
9 Accordingly, accurate phenotyping of HF patients usually requires the incorporation of  
10  
11 data that crosses clinical encounters. In addition, although symptoms are an essential  
12  
13 reflection of clinical status, they are poorly captured in structured data. Suboptimal  
14  
15 recognition of comorbidities like valvular heart disease can also impact disease  
16  
17 trajectory and risk for future cardiovascular events.  
18  
19  
20  
21  
22

23 Our findings represent an important advance for RWE studies that include HF patients.  
24  
25 Notably, the only way to ascertain comparative accuracy between data sources and  
26  
27 technologies in a domain is to test it. Accuracy consists of both recall and precision, and  
28  
29 in the case of many health conditions, recall can fall below 50% when one relies solely  
30  
31 upon the problem list.<sup>25,26</sup>  
32  
33  
34

35 In the current study, we were able to focus on both precision and recall through use of  
36  
37 the  $F_1$  score. Despite availability of SNOMED codes for HF<sub>r</sub>EF and HF<sub>p</sub>EF, along with a  
38  
39 similar code for HF<sub>mr</sub>EF, such codes were rarely included. Documentation of a HF code  
40  
41 using structured data was only found 46.9% of the time when there was clear evidence  
42  
43 of HF in the chart. We postulate that the low accuracy of structured data for disease  
44  
45 subtypes at least partially relates to how the data is likely to be used. A physician may  
46  
47 look within notes to understand HF subtype. Information entered into problem lists and  
48  
49 claims may be more to provide a high-level understanding of disease burden. Granular  
50  
51 billing codes may be a low priority for physicians if claims are reimbursed with the non-  
52  
53  
54  
55  
56  
57

1  
2  
3 granular HF code. Furthermore, because addition of diagnoses to the problem list is not  
4  
5 a requirement, the problem list may not be specific or updated. This contrasts with  
6  
7 clinical notes, where detailed documentation is usually performed to communicate a  
8  
9 care plan and is a medical-legal requirement.  
10  
11

12  
13 When low-accuracy and non-granular data are utilized, there are several potential  
14  
15 consequences. Missingness can result in selection bias, particularly if sicker patients  
16  
17 have more frequent encounters, higher rates of specialty care, and more complete  
18  
19 documentation. Depending on the study question, use of structured data alone to  
20  
21 identify certain subgroups may be inadvisable, since these data have a low recall for  
22  
23 specific clinical concepts such as ST-elevation myocardial infarction and HFrEF.<sup>27</sup> Even  
24  
25 advanced approaches (e.g., NLP) may result in poor accuracy, as illustrated in this  
26  
27 study, where HFpEF required AI-based inference for proper identification. Collectively,  
28  
29 this highlights that not all data sources and technologies are the same; therefore,  
30  
31 accuracy testing may be required for rigorous RWE generation. Furthermore, given the  
32  
33 growth in RWE to support new drug indications, post-marketing surveillance, and  
34  
35 decision-making regarding reimbursement, such inaccuracies may have a profound  
36  
37 impact on large numbers of patients.  
38  
39  
40  
41  
42  
43

44 Even though standard dictionaries and clinical terms related to cardiovascular medicine  
45  
46 were used, there is a need to test the two analytic methods using different EHRs across  
47  
48 a broader set of community and referral practices. With numerous EHRs available and  
49  
50 practitioner-to-practitioner variability in documentation accuracy, efforts like the one  
51  
52 described here represent an important means of strengthening data quality.  
53  
54  
55  
56  
57



1  
2  
3 Importantly, this study has several limitations. First, we used data from a single health  
4 system, with results that may not be generalizable to other populations. Second, the  
5 study protocol required the selection of patients enriched with cardiovascular disease to  
6 make the study feasible, with manual chart abstraction conducted to ensure the  
7 accuracy of results. While selection criteria were applied to both structured and  
8 unstructured data, it is possible that this could have biased results in a way that favored  
9 structured data since a larger proportion of patients with HF on the problem list may  
10 have been included than if the sample had been created randomly. In addition, the  
11 specific filters used likely led to a higher-than-expected proportion of HF<sub>r</sub>EF patients  
12 (compared to those with HF<sub>mr</sub>EF and HF<sub>p</sub>EF). Second, the study required laborious  
13 manual annotation of thousands of records. Such a sample size is adequate for high-  
14 prevalence conditions, but would likely require adjustment for low-prevalence conditions  
15 with low concept occurrence rates. Finally, the study did not include clinical outcome  
16 assessment; rather, it was designed to compare data sources and processing methods.  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

## 36 **Conclusion**

37  
38  
39  
40 As RWE is increasingly used to analyze patient subgroups, inform clinical decision-  
41 making, and influence regulatory and reimbursement decisions, data reliability and  
42 evidence validity are of critical importance. Use of a traditional approach was associated  
43 with low data accuracy. While much greater accuracy was observed with AI-based  
44 methods, it depended upon the technology utilized. These findings highlight the  
45 importance of using data fit-for-purpose to the research question posed. In addition,  
46 they suggest that accuracy testing should be part of any EHR-based study that includes  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

1  
2  
3 HF patients. Finally, unstructured data and a technology-based approach to data  
4  
5 extraction may be required in some studies to achieve sufficient accuracy, depending  
6  
7 upon the clinical assertion being tested.  
8  
9

## 10 11 **Acknowledgments** 12

13  
14 We are grateful for comments from Jacob Abraham, MD, Medical Director at  
15  
16 Providence Heart Institute's Center for Advanced Heart Disease, and Yuri Quintana,  
17  
18 MD, Chief of, Division of Clinical Informatics at the Beth Israel Deaconess Medical  
19  
20 Center. Editorial support was provided by Liam Gillies, Ph.D., CMPP, of Cactus Life  
21  
22 Sciences (part of Cactus Communications), funded by Amgen Inc.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

## Contributors

ARG and DR drafted the manuscript. ARG, KLM, RED, DR, and TJG critically reviewed the manuscript. ARG, RED, DR, and TJG provided clinical insight.

## Funding

A research grant supported this work from Amgen Inc. DR was partly supported by the US Food and Drug Administration (FDA) under Award Number IIP-2024958 and the National Center for Advancing Translational Sciences of the NIH under Award Number R44TR002437. The content is solely the responsibility of the authors and does not necessarily represent the official views of Amgen, the FDA, or the NIH.

## Competing interests

KLM and RED are employees and stockholders of Amgen Inc. DR is an employee and stockholder of Verantos, Inc. ARG has received research support from Abbott and TJG has no competing interests to declare.

## Ethics Approval

This study has been independently reviewed and accepted for exemption in accordance with 45 CFR 46.101(b)(4ii).

## Provenance and peer review

Not commissioned, externally peer reviewed.

## Data sharing statement

No additional data are available.

## Supplemental Materials

Supplemental Methods

## REFERENCES

1. Thomas H, Diamond J, Vieco A, Chaudhuri S, Shinnar E, Cromer S, Perel P, Mensah GA, Narula J, Johnson CO, et al. Global atlas of cardiovascular disease 2000-2016: the path to prevention and control. *Glob Heart*. 2018;13:143-163. doi: 10.1016/j.gheart.2018.09.511
2. Nichols M, Townsend N, Scarborough P, Rayner M. Cardiovascular disease in Europe 2014: epidemiological update. *Eur Heart J*. 2014;35:2950-2959. doi: 10.1093/eurheartj/ehu299
3. McMurray JJ, Packer M, Desai AS, Gong J, Lefkowitz MP, Rizkala AR, Rouleau JL, Shi VC, Solomon SD, Swedberg K, et al. Angiotensin-neprilysin inhibition versus enalapril in heart failure. *N Engl J Med*. 2014;371:993-1004. doi: 10.1056/NEJMoa1409077
4. McMurray JJV, Solomon SD, Inzucchi SE, Køber L, Kosiborod MN, Martinez FA, Ponikowski P, Sabatine MS, Anand IS, Bělohávek J, et al. Dapagliflozin in patients with heart failure and reduced ejection fraction. *N Engl J Med*. 2019;381:1995-2008. doi: 10.1056/NEJMoa1911303
5. Packer M, Anker SD, Butler J, Filippatos G, Pocock SJ, Carson P, Januzzi J, Verma S, Tsutsui H, Brueckmann M, et al. Cardiovascular and renal outcomes with empagliflozin in heart failure. *N Engl J Med*. 2020;383:1413-1424. doi: 10.1056/NEJMoa2022190.
6. Swedberg K, Komajda M, Böhm M, Borer JS, Ford I, Dubost-Brama A, Lerebours G, Tavazzi L, SHIFT Investigators. Ivabradine and outcomes in chronic heart

- 1  
2  
3 failure (SHIFT): a randomised placebo-controlled study. *Lancet*. 2010;376:875-  
4 885. doi: 10.1016/S0140-6736(10)61198-1  
5  
6  
7  
8 7. Stone GW, Lindenfeld J, Abraham WT, Kar S, Lim DS, Mishell JM, Whisenant B,  
9 Grayburn PA, Rinaldi M, Kapadia SR, et al. Transcatheter mitral-valve repair in  
10 patients with heart failure. *N Engl J Med*. 2018;379:2307-2318. doi:  
11 10.1056/NEJMoa1806640  
12  
13  
14  
15  
16  
17 8. Shah KS, Xu H, Matsouka RA, Bhatt DL, Heidenreich PA, Hernandez AF,  
18 Devore AD, Yancy CW, Fonarow GC. Heart failure with preserved, borderline,  
19 and reduced ejection fraction: 5-year outcomes. *J Am Coll Cardiol*.  
20 2017;70:2476-2486.  
21  
22  
23  
24  
25  
26 9. H.R.34 - 21st Century Cures Act of 2016. Public Law No. 114-255. Section 3022.  
27 Available at: <https://www.congress.gov/bill/114th-congress/house-bill/34/>.  
28 Accessed May 4, 2020.  
29  
30  
31  
32  
33 10. Pulini AA, Caetano GM, Clautiaux H, Vergeron L, Pitts PJ, Katz G. Impact of  
34 real-world data on market authorization, reimbursement decision & price  
35 negotiation [published online ahead of print August 28, 2020]. *Ther Innov Regul*  
36 *Sci*. 2020. doi: 10.1007/s43441-020-00208-1  
37  
38  
39  
40  
41  
42 11. Hernandez-Boussard T, Monda KL, Crespo BC, Riskin D. Real world evidence in  
43 cardiovascular medicine: ensuring data validity in electronic health record-based  
44 studies. *J Am Med Inform Assoc*. 2019;26:1189-1194. doi:  
45 10.1093/jamia/ocz119.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 12. McCormick N, Lacaille D, Bhole V, Avina-Zubieta JA. Validity of heart failure  
4 diagnoses in administrative databases: a systematic review and meta-analysis.  
5  
6 *PLoS One*. 2014;9:e104519. doi: 10.1371/journal.pone.0104519.  
7  
8  
9
- 10 13. Alqaisi F, Williams LK, Peterson EL, Lanfear DE. Comparing methods for  
11 identifying patients with heart failure using electronic data sources. *BMC Health*  
12 *Serv Res*. 2009;9:237. doi: 10.1186/1472-6963-9-237  
13  
14  
15  
16
- 17 14. Xu Y, Lee S, Martin E, D'souza AG, Doktorchik CTA, Jiang J, Lee S, Eastwood  
18 CA, Fine N, Hemmelgarn B, et al. Enhancing ICD-code-based case definition for  
19 heart failure using electronic medical record data. *J Card Fail*. 2020;26:610-617.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32
- 33 15. [https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world-data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory)  
34 [world-data-assessing-electronic-health-records-and-medical-claims-data-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world-data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory)  
35 [support-regulatory](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world-data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory)  
36  
37  
38  
39
- 40 16. Cunningham H, Tablan V, Roberts A, Bontcheva K. Getting more out of  
41 biomedical documents with GATE's full lifecycle open source text analytics. *PLoS*  
42 *Comput Biol*. 2013;9(2):e1002854.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57
- 58 17. Van Rijsbergen CJ. Information Retrieval. 2nd ed. Butterworth-Heinemann. 1979.  
59  
60
- 61 18. Bozkurt B, Coats AJ, Tsutsui H, et al. Universal Definition and Classification of  
62 Heart Failure: A Report of the Heart Failure Society of America, Heart Failure  
63 Association of the European Society of Cardiology, Japanese Heart Failure  
64 Society and Writing Committee of the Universal Definition of Heart Failure. *J*  
65 *Card Fail*. 2021 Mar 1:S1071-9164(21)00050-6. doi:  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100

- 1  
2  
3 19. Bielinski SJ, Pathak J, Carrell DS, et al. A Robust e-Epidemiology Tool in  
4  
5 Phenotyping Heart Failure with Differentiation for Preserved and Reduced  
6  
7 Ejection Fraction: the Electronic Medical Records and Genomics (eMERGE)  
8  
9 Network. *J Cardiovasc Transl Res*. 2015 Nov;8(8):475-83.  
10  
11  
12 20. Blecker S, Katz SD, Horwitz LI, Kuperman G, Park H, Gold A, Sontag D.  
13  
14 Comparison of approaches for heart failure case identification from electronic  
15  
16 health record data. *JAMA Cardiol*. 2016;1:1014-1020. doi:  
17  
18 10.1001/jamacardio.2016.3236  
19  
20  
21 21. Ng K, Steinhubl SR, deFilippi C, Dey S, Stewart WF. Early detection of heart  
22  
23 failure using electronic health records: practical implications for time before  
24  
25 diagnosis, data diversity, data quantity, and data density. *Circ Cardiovasc Qual*  
26  
27 *Outcomes*. 2016;9:649-658. doi: 10.1161/CIRCOUTCOMES.116.002797  
28  
29  
30  
31 22. Tison GH, Chamberlain AM, Pletcher MJ, Dunlay SM, Weston SA, Killian JM,  
32  
33 Olgin JE, Roger VL. Identifying heart failure using EMR-based algorithms. *Int J*  
34  
35 *Med Inform*. 2018;120:1-7. doi: 10.1016/j.ijmedinf.2018.09.016  
36  
37  
38 23. U.S. Food and Drug Administration. Framework for FDA's Real-World Evidence  
39  
40 Program. 2018. Available at: <https://www.fda.gov/media/120060/download>.  
41  
42 Accessed July 26, 2020.  
43  
44  
45 24. Parsons A, McCullough C, Wang J, Shih S. Validity of electronic health record-  
46  
47 derived quality measurement for performance monitoring. *J Am Med Inform*  
48  
49 *Assoc*. 2012;19:604-609. doi: 10.1136/amiainl-2011-000557  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



- 1  
2  
3 25. Luna D, Franco M, Plaza C, Otero C, Wassermann S, Gambarte ML, Giunta D,  
4  
5 de Quirós FGB. Accuracy of an electronic problem list from primary care  
6  
7 providers and specialists. *Stud Health Technol Inform*. 2013;192:417-421.  
8  
9
- 10 26. Singer A, Yakubovich S, Kroeker AL, Dufault B, Duarte R, Katz A. Data quality of  
11  
12 electronic medical records in Manitoba: do problem lists accurately reflect chronic  
13  
14 disease billing diagnoses? *J Am Med Inform Assoc*. 2016;23:1107-1112. doi:  
15  
16 10.1093/jamia/ocw013  
17  
18
- 19 27. Makam AN, Lanham HJ, Batchelor K, Samal L, Moran B, Howell-Stampley T,  
20  
21 Kirk L, Cherukuri M, Santini N, Leykum LK, et al. Use and satisfaction with key  
22  
23 functions of a common commercial electronic health record: a survey of primary  
24  
25 care providers. *BMC Med Inform Decis Mak*. 2013;13:86. doi: 10.1186/1472-  
26  
27 6947-13-86  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

Table 1. Prespecified heart failure–specific concepts extracted from the electronic health record.

High Priority Conditions	Comorbidities	Symptoms	Findings	Medications
Congestive HF	Myocardial infarction	Angina	LVEF	Carvedilol
HF with reduced EF	Atrial fibrillation	Chest pain		Lisinopril
HF with mid-range EF	Aortic regurgitation	Dyspnea		Metoprolol
HF with preserved EF	Mitral regurgitation	Fatigue		Furosemide
	Tricuspid regurgitation	Palpitations		

HF, heart failure; EF, ejection fraction; LVEF, left ventricular ejection fraction.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

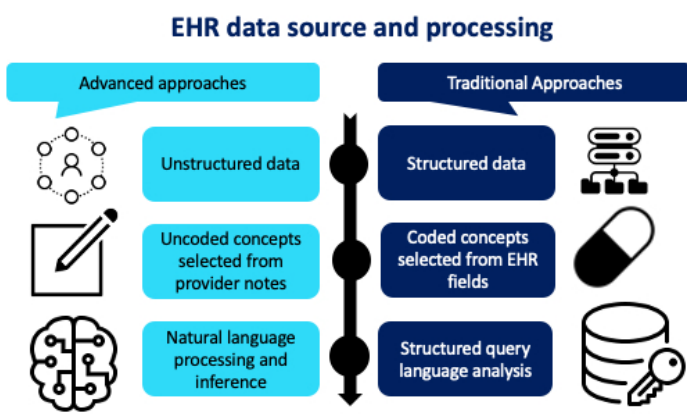


Figure 1. Comparison of traditional and advanced real-world evidence approaches. EHR, electronic health record.

338x190mm (54 x 54 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

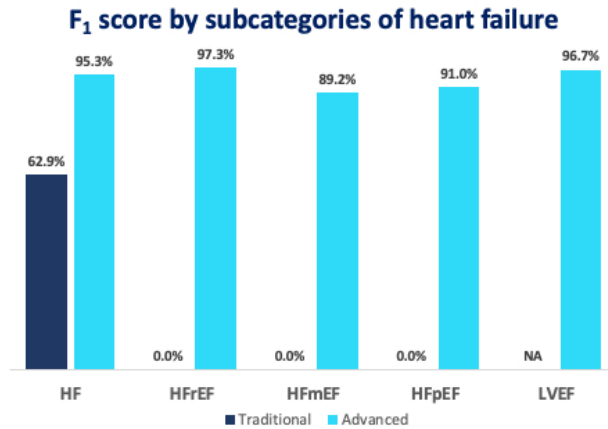


Figure 2. F1 scores for heart failure diagnoses. \*F1-score could not be calculated due to lack of data for precision. †Structured data recall is not applicable for ejection fraction because no code was available within the problem list. HF, heart failure; HFmEF, heart failure with mildly-reduced ejection fraction; HFpEF, heart failure with preserved ejection fraction; HFrEF, heart failure with reduced ejection fraction; LVEF, left ventricular ejection fraction; 0% reflects a measured value and indicates the availability of the diagnosis code in the EHR dropdown versus N/A, not applicable, which refers to a diagnosis without available code in the relevant codeset.

338x190mm (54 x 54 DPI)

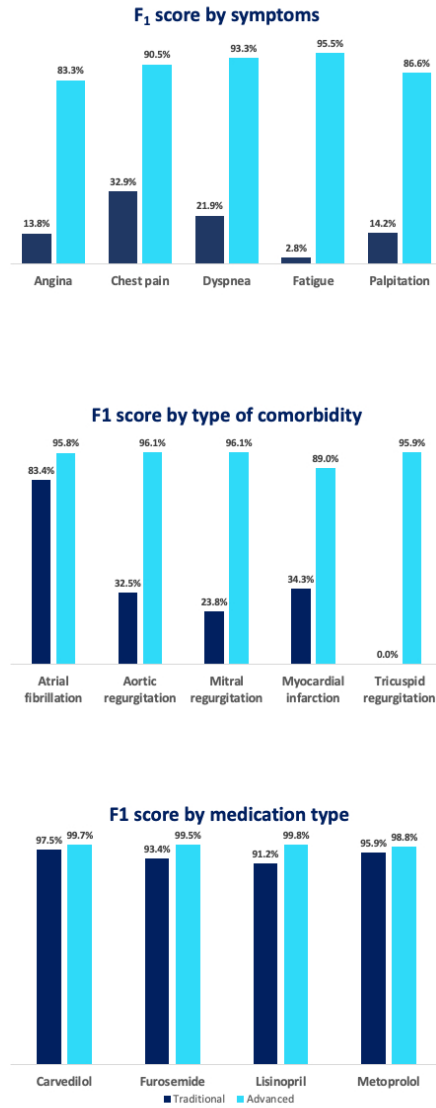


Figure 3. F1 scores for (A) symptoms, (B) medications, and (C) comorbid conditions. \*F1 score could not be calculated due to a lack of data for precision. N/A, not applicable.

254x428mm (72 x 72 DPI)

Supplementary Table 1. Cohort identification of heart failure diagnoses, left ventricular ejection fraction, heart failure medications, symptoms, and comorbid cardiovascular conditions

	Traditional approach			Advanced approach			Concept occurrence	Encounter occurrence	P-value
	Recall, %	Precision, %	F <sub>1</sub> score, %	Recall, %	Precision, %	F <sub>1</sub> score, %			
<b>HF diagnosis</b>									
HF	46.9	95.4	62.9	96.0	94.7	95.3	265	155	<0.001
HFrEF	0	N/A*	N/A <sup>†</sup>	94.8	100.0	97.3	382	124	N/A <sup>§</sup>
HFmrEF	0	N/A*	N/A <sup>†</sup>	80.4	100.0	89.2	62	35	N/A <sup>§</sup>
HFpEF	0	N/A*	N/A <sup>†</sup>	83.5	100.0	91.0	103	55	N/A <sup>§</sup>
<b>LVEF</b>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	93.7	100.0	96.7	677	238	N/A <sup>§</sup>
<b>HF medications</b>									
Carvedilol	95.1	100.0	97.5	99.7	99.7	99.7	407	141	<0.001
Furosemide	87.7	100.0	93.4	99.3	99.8	99.5	1572	371	0.116
Lisinopril	83.9	100.0	91.2	99.7	99.9	99.8	1068	386	<0.001
Metoprolol	92.2	100.0	95.9	97.7	100.0	98.8	1370	397	<0.001
<b>Symptoms</b>									

Angina	7.8	60.0	13.8	84.4	82.3	83.3	265	155	<0.00 1
Chest pain	21.4	70.8	32.9	95.4	86.1	90.5	2332	756	<0.00 1
Dyspnea	12.7	78.2	21.9	94.7	92.0	93.3	4474	832	<0.00 1
Fatigue	1.4	75.0	2.8	96.5	94.5	95.5	1711	371	<0.00 1
Palpitation	8.2	52.9	14.2	90.9	82.6	86.6	896	493	<0.00 1
<b>Comorbid cardiovascular conditions</b>									
Atrial fibrillation	72.2	98.7	83.4	93.0	98.7	95.8	1214	222	<0.00 1
Aortic regurgitation	19.4	100.0	32.5	92.5	100.0	96.1	153	90	<0.00 1
Mitral regurgitation	13.5	97.1	23.8	92.8	99.6	96.1	483	185	<0.00 1
Myocardial infarction	21.1	90.9	34.3	95.5	83.4	89.0	1220	578	<0.00 1
Tricuspid regurgitation	0	N/A*	N/A†	92.2	100.0	95.9	162	78	N/A§

\*These elements did not occur when using the traditional approach. †F<sub>1</sub> scores could not be calculated due to a lack of data for precision. ‡Structured data recall is not applicable for ejection fraction because there was no code available within the problem list. §P-value could not be calculated due to the unavailability of F<sub>1</sub> scores for the traditional approach. P-values are derived from the chi-square test.

1  
2  
3 HF, heart failure; HFmrEF, heart failure with mildly reduced ejection fraction; HFpEF, heart failure with  
4 preserved ejection fraction; HFrEF, heart failure with reduced ejection fraction; LVEF, left ventricular  
5 ejection fraction; N/A, not applicable.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## SUPPLEMENTAL MATERIAL

### Reference Standard

Traditional and advanced approaches were tested against a reference standard for physician encounters. The reference standard consisted of an independent review, with manual annotation of relevant HF-specific features, including 19 unique HF-specific concepts. For each encounter, two independent clinical annotators labeled each concept and all metadata for that concept. For example, an annotator might mark the text "DOE over last month" as dyspnea on exertion, experienced = true, current = true, relative date = 1 month. Concept occurrence was defined as the sum of all concept occurrences, allowing for multiple occurrences per encounter. Encounter occurrence was defined as the number of encounters with at least one occurrence of the concept.

Given that many concepts, such as LVEF are specific to a point in time, concepts were tested at the encounter level. For example, if a patient had an LVEF of 30% in an encounter, the data extraction would only be annotated as correct if it identified "LVEF 30%" in that specific encounter. This reference standard was used to determine accuracy of automated extracted data and structured data. Specifically, this reference standard was used to calculate recall and precision for these individual features for traditional and advanced approaches.

# BMJ Open

## A Comparison of Traditional and Artificial-Intelligence Based Heart Failure Phenotyping to Enable Real-World Evidence

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2023-073178.R1
Article Type:	Original research
Date Submitted by the Author:	29-Jun-2023
Complete List of Authors:	Garan, Arthur Reshad; Harvard Medical School Monda , Keri; Amgen Inc Dent-Acosta, Ricardo ; Amgen Inc Riskin , Daniel; Verantos Gluckman, Ty; Providence St Joseph Health
<b>Primary Subject Heading</b>:	Health informatics
Secondary Subject Heading:	Cardiovascular medicine
Keywords:	Heart failure < CARDIOLOGY, Cardiac Epidemiology < CARDIOLOGY, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

1  
2  
3 **A Comparison of Traditional and Artificial-Intelligence Based Heart Failure Phenotyping**  
4 **to Enable Real-World Evidence**  
5  
6  
7

8  
9 **Short title:** Real-World Evidence in Heart Failure  
10

11  
12 **Authors:** A. Reshad Garan, MD, MS; Keri L. Monda, Ph.D.; Ricardo E. Dent-Acosta; MD;  
13  
14 Dan Riskin, MD; Ty J. Gluckman, MD, MHA  
15  
16

17  
18 **Affiliations:** Department of Medicine, Division of Cardiology, Beth Israel Deaconess Medical  
19  
20 Center, Harvard Medical School, Boston, MA, USA (A.R.G.). The Center for Observational  
21  
22 Research and Medical Affairs, Amgen Inc., Thousand Oaks, CA, USA (K.L.M., R.E.D.-A.).  
23  
24 Department of Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, NC,  
25  
26 USA (K.L.M.). Verantos, Inc., Menlo Park, CA, USA (D.R.). Department of Surgery, Stanford  
27  
28 University School of Medicine, Stanford, CA, USA (D.R.). Center for Cardiovascular Analytics,  
29  
30 Research and Data Science (CARDS), Providence Heart Institute, Providence St. Joseph Health,  
31  
32 Portland, OR, USA (T.J.G.)  
33  
34  
35  
36

37  
38 **Address for correspondence:** A. Reshad Garan, MD, MS, Department of Medicine, Division of  
39  
40 Cardiology, Beth Israel Deaconess Medical Center, 185 Pilgrim Road Deaconess 3, Boston, MA  
41  
42 02215, USA. Phone: 617/632-7737. Fax: 617/632-7620. Email: [agaran@bidmc.harvard.edu](mailto:agaran@bidmc.harvard.edu)  
43  
44  
45

46 **Word count:** 2673 words  
47  
48

49 **Keywords:** artificial intelligence; heart failure; phenotype; real-world evidence; electronic health  
50  
51 record.  
52  
53

54  
55 **Abstract**  
56  
57

1  
2  
3 **Objective:** Quantitatively evaluate the quality of data underlying real-world evidence (RWE) in  
4 heart failure (HF).  
5  
6

7 **Design:** Retrospective comparison of accuracy in identifying HF patients and phenotypic  
8 information was made using traditional (i.e., structured query language applied to structured  
9 EHR data) and advanced (i.e., AI applied to unstructured EHR data) RWE approaches. The  
10 performance of each approach was measured by the harmonic mean of precision and recall (F1  
11 score) using manual annotation of medical records as a reference standard.  
12  
13  
14  
15  
16  
17  
18

19 **Setting:** EHR data from a large academic healthcare system in North America between 2015 and  
20 2019, with an expected catchment of approximately 500,000 patients.  
21  
22  
23

24 **Population:** 4288 encounters for 1155 patients aged 18 to 85 years, with 472 patients identified  
25 as having HF.  
26  
27

28 **Outcome measures:** HF and associated concepts, such as comorbidities, left ventricular ejection  
29 fraction, and selected medications.  
30  
31  
32  
33

34 **Results:** The average F1 scores across 19 HF-specific concepts were 49.0% and 94.1% for the  
35 traditional and advanced approaches, respectively ( $P < 0.001$  for all concepts with available data).  
36 The absolute difference in F<sub>1</sub> score between approaches was 45.1% (98.1% relative increase in  
37 F<sub>1</sub> score using the advanced approach). The advanced approach achieved superior F1 scores for  
38 HF presence, phenotype, and associated comorbidities. Some phenotypes, such as HFpEF,  
39 revealed dramatic differences in extraction accuracy based on technology applied, with a 4.9%  
40 F<sub>1</sub> score when using natural language processing (NLP) alone and a 91.0% F<sub>1</sub> score when using  
41 NLP plus AI-based inference.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **Conclusions:** A traditional RWE generation approach resulted in low data quality in HF patients.  
4  
5 While an advanced approach demonstrated high accuracy, the results varied dramatically based  
6  
7 on extraction techniques. For future studies, advanced approaches and accuracy measurement  
8  
9 may be required to ensure data are fit-for-purpose.  
10  
11  
12

### 13 **Strengths and limitations of this study**

- 14  
15  
16  
17 • Using RWE for HF patients requires demonstrating that the data source and technologies  
18  
19 result in accurate data.
- 20  
21  
22 • Natural language processing alone lacked context from the longitudinal record, limiting  
23  
24 phenotype identification and study validity.
- 25  
26  
27 • Findings suggest that advanced methods can enable high-validity RWE for heart failure  
28  
29 patients.
- 30  
31 • The use of data from a single healthcare system may limit generalizability to other  
32  
33 populations.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

## INTRODUCTION

Heart failure (HF) is a major public health problem with significant associated morbidity, mortality, and cost.<sup>1,2</sup> Despite the availability of novel drugs and devices, morbidity and mortality in HF rivals many malignancies, with a 5-year survival rate as low as 50%.<sup>3-8</sup>

Randomized controlled trials (RCTs) have traditionally been used to assess the safety and efficacy of new therapies and represent a cornerstone for regulatory approval. However, RCTs are frequently conducted in highly selected populations, typically younger, healthier, and less diverse than patients treated in clinical practice. Furthermore, such trials often include patients with an established HF diagnosis, receiving guideline-directed medical therapy at tertiary centers, and may not represent the broader HF population. Because HF is a clinically heterogeneous syndrome with numerous etiologies and phenotypes, studying this population can be particularly difficult.

Real-world evidence (RWE) has held promise as a potential means to assess therapeutic benefit outside of clinical trials, with sufficient power to characterize therapeutic impact in HF subgroups. Accordingly, RWE can complement RCTs, extending the findings to patient populations that may have been excluded from or insufficiently enrolled in pivotal trials. To accelerate these and similar precision medicine goals, the 21st Century Cures Act was passed in 2016, which required the United States Food and Drug Administration (FDA) to develop guidance supporting the use of RWE in new drug indications and post-marketing surveillance.<sup>9</sup> In addition, payors have increasingly utilized RWE to inform reimbursement decisions and are increasingly demanding credible evidence.<sup>10</sup>

1  
2  
3 Not surprisingly, the quality of RWE hinges on how well real-world data are collected,  
4 processed<sup>11</sup>, and used to inform study questions. Such is the case in HF, where accurate  
5 identification of patients in administrative and other structured data sets is an ongoing focus.<sup>12-14</sup>  
6  
7 Traditional methods of identifying HF patients rely on querying diagnosis codes and structured  
8 data in the electronic health record (EHR) or medical claims. Conversely, artificial intelligence  
9 (AI) applied to unstructured data represents a novel method of analyzing the medical record.  
10  
11 Because of the importance of data reliability in RWE and the potential to use unstructured data to  
12 achieve data enrichment<sup>15</sup>, we sought to compare the accuracy achieved by traditional RWE  
13 methods versus advanced AI approaches in identifying a range of HF-specific data elements  
14 from the medical record.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

## 26 27 **METHODS**

28  
29  
30 The study design is outlined in Figure 1. Varied data sources and applied technologies were used  
31 to assess data reliability in patients with risk factors for HF. Leveraging manual chart abstraction  
32 as the reference standard, comparisons were made between the two methods. The first method  
33 used structured EHR data (e.g., diagnosis codes and problem lists) and standard query  
34 techniques, defined as the 'traditional approach'. The second used unstructured EHR data (e.g.,  
35 narratives from primary care and specialty notes) and AI techniques, described as the 'advanced  
36 approach' (Figure 1). The primary objective was measurement of the accuracy of identified HF-  
37 specific elements using traditional and advanced approaches. We hypothesized that the advanced  
38 approach would better identify key HF-specific elements than the traditional approach. Data  
39 were deidentified before study initiation, and the study was determined not to be human subjects  
40 research. Both natural language processing (NLP) and machine-learned inference technologies  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57



1  
2  
3 used in the advanced approach were provided by Verantos, Inc. (Menlo Park, CA, USA). The  
4  
5 core of AI is a deterministic NLP layer. This layer is built on top of the GATE NLP  
6  
7 architecture.<sup>16</sup> The architecture is used to construct a flexible pipeline for processing incoming  
8  
9 text against English language syntactical rules augmented with a lexicon based on a clinical  
10  
11 vocabulary. The AI-based inference was applied during data processing. Millions of machine-  
12  
13 learned and manually curated associations enable disambiguation and identification of clinically  
14  
15 relevant concepts. As an example of AI-based inference, a patient with HF on the problem list  
16  
17 and a narrative encounter describing “EF 60%” would not be interpreted by NLP as having HF  
18  
19 with preserved ejection fraction (HFpEF) since the text does not have sufficient information to  
20  
21 identify this condition. On the other hand, AI-based inference would infer HFpEF based on  
22  
23 disparate information in the record.  
24  
25  
26  
27  
28

### 29 **EHR Data Source and Processing**

30  
31  
32  
33 EHR data from primary care encounters between 2011 and 2018 were deidentified and securely  
34  
35 transferred to a cloud-based server for analysis. The data set consisted of both structured data  
36  
37 (e.g., medical conditions, procedures performed, medications, and problem lists) and  
38  
39 unstructured data (e.g., narrative notes from primary care providers and specialists, telephone  
40  
41 visits, and other narrative text) (Figure 2).  
42  
43  
44

45  
46 As the study aimed to test the accuracy of different RWE approaches and not treatment  
47  
48 effectiveness, the cohort was enriched for patients with suspected HF based on comorbidities and  
49  
50 medications. Specifically, the following filters were applied: records containing both narrative  
51  
52 and structured components; narrative length 1,000 characters or more; and at least one of the  
53  
54  
55  
56  
57

1  
2  
3 following problems or medications in structured or unstructured data: myocardial infarction,  
4 congestive heart failure, or carvedilol (Figure 1).  
5  
6  
7

8  
9 A prespecified set of clinical concepts pertinent to patients with HF was extracted using  
10 traditional and advanced techniques (Table 1). Problem lists were mapped to Systematized  
11 Nomenclature of Medicine (SNOMED) ontology, and unadjudicated claims were mapped to  
12 ICD-10 codes. Standard sets of individual codes were used to represent each concept. With the  
13 advanced approach, inference incorporating pattern recognition was utilized to identify  
14 potentially missing or ignored concepts within the text (e.g., HF being likely in patients with  
15 dyspnea and pitting edema on a diuretic). Specifically, no narrative coding took place before the  
16 AI algorithm was used; instead, it was applied directly to the narrative text and then mapped by  
17 the algorithm to the SNOMED ontology. Next, manual chart abstraction using the same  
18 SNOMED code set was used as a reference to assess the accuracy of the coding by the AI  
19 algorithm. Engineers were blinded to validation data and its corresponding chart abstraction.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34

### 35 **Study End Points and Statistical Analysis**

36  
37  
38 The primary endpoint was the  $F_1$  score for traditional and advanced approaches. The  $F_1$  score is  
39 an accuracy measure that combines recall and precision; more specifically, it is the weighted  
40 harmonic mean of these two measures. Secondary endpoints were recall (i.e., the proportion of  
41 patients correctly identified as having the condition, akin to sensitivity) and precision (i.e., the  
42 proportion of patients with HF and its subtypes correctly identified divided by the total number  
43 of patients identified in each cohort akin to positive predictive value)<sup>17,18</sup> for the traditional and  
44 advanced approaches. The reference standard used to evaluate accuracy of the traditional and  
45 advanced approaches was manual chart abstraction. For each encounter, two independent clinical  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

1  
2  
3 annotators labeled each concept and all metadata for that concept. Annotators were blinded to  
4  
5 each other's annotations, and inter-rater agreement was measured by Cohen's kappa score.  
6  
7  
8 Further description of the reference standard methodology is provided in the Supplemental  
9  
10 Material. Results were summarized using descriptive statistics, and percentages were calculated  
11  
12 for categorical variables. Differences in  $F_1$  scores between traditional and advanced approaches  
13  
14 were analyzed using the chi-square test; associated  $P$ -values were reported.  
15  
16

### 17 18 **Patient and Public Involvement**

19  
20  
21 Data were deidentified before study initiation, and the study was determined not to be human  
22  
23 subjects research. As a result, no patients were recruited for study participation. The research  
24  
25 question and study goal of highlighting methods for improving RWE use were driven by  
26  
27 recognition that improvements in use of RWE to inform new drug indications, post-marketing  
28  
29 surveillance, and reimbursement decisions would ultimately result in patient benefit.  
30  
31  
32

### 33 34 **RESULTS**

35  
36  
37 A total of 4288 encounters for 1155 patients were examined, of which 472 patients with HF were  
38  
39 identified. Of these, 382 had HF with reduced ejection fraction (HF<sub>r</sub>EF), 35 had HF with mildly  
40  
41 reduced ejection fraction (HF<sub>mr</sub>EF), and 55 had HF with preserved ejection fraction (HF<sub>p</sub>EF).  
42  
43  
44 The reference standard Cohen's kappa score was 0.95, suggesting high validity.  
45  
46

47  
48 Supplementary Table 1 reports the  $F_1$  score, recall, and precision results achieved with both  
49  
50 approaches. Figure 3 graphically presents  $F_1$  scores for HF diagnoses and Figure 4 includes  $F_1$   
51  
52 scores for symptoms, medications, and comorbid conditions. Overall, accuracy was significantly  
53  
54 greater for the advanced approach (AI applied to unstructured EHR data) than for the traditional  
55  
56  
57

1  
2  
3 approach (structured query language applied to structured EHR data) (Supplementary Table 1;  
4  
5 Figure 3; Figure 4), with an absolute difference of 45.1%.

6  
7  
8  
9 With the traditional approach, recall for any HF diagnosis was 46.9% (i.e., 53.1% of patients  
10  
11 with HF were missed entirely) and precision was 95.4%, resulting in an  $F_1$  score of 62.9%  
12  
13 ( $P<0.001$ ). In contrast, with the advanced approach, recall for any HF diagnosis was 96.0% and  
14  
15 precision was 94.7%, resulting in an  $F_1$ -score of 95.3% ( $P<0.001$  when  $F_1$  scores for the two  
16  
17 approaches were compared) (Supplementary Table 1; Figure 3). Among HF phenotypes, recall  
18  
19 with the advanced approach was highest with HFrEF, followed by HFpEF and HFmrEF;  
20  
21 precision was 100% for all phenotypes. With the traditional approach,  $F_1$  scores could not be  
22  
23 calculated for HFrEF, HFmrEF, and HFpEF because only less granular HF codes were used  
24  
25 (Supplementary Table 1).

26  
27  
28  
29  
30  
31 Accuracy in identifying left ventricular ejection fraction (LVEF) was similarly high with the  
32  
33 advanced approach, with an  $F_1$  score of 96.7%. Data could not be extracted for LVEF with the  
34  
35 traditional approach because no such codes were available within the EHR, nor did a mechanism  
36  
37 to encode LVEF within the problem list or unadjudicated claims exist (Supplementary Table 1;  
38  
39 Figure 3).

40  
41  
42  
43 Accurate identification of HF symptoms was greater with the advanced approach ( $P<0.001$ )  
44  
45 (Supplementary Table 1; Figure 4A). Whereas identification of commonly prescribed HF  
46  
47 medications was high with both approaches (Supplementary Table 1; Figure 4B), identification  
48  
49 of cardiovascular comorbidities was higher in all cases with the advanced approach ( $P<0.001$ )  
50  
51 (Supplementary Table 1; Figure 4C).

1  
2  
3 Data concept extraction with the advanced approach greatly depended upon the technology used.  
4  
5 For example, NLP, which ends at the sentence boundary, was only able to identify HFpEF with  
6  
7 an F<sub>1</sub> score of 4.9% because "HFpEF" or "heart failure with preserved ejection fraction" was  
8  
9 rarely written. Conversely, inference, which can find related items from the longitudinal record,  
10  
11 was able to identify both "HF" and "normal ejection fraction" as separate annotations for HFpEF  
12  
13 with an F<sub>1</sub>-score of 91.0% (Supplementary Table 1; Figure 3).  
14  
15  
16  
17

## 18 **DISCUSSION**

19  
20  
21 The utilization of RWE has grown substantially in recent years, driven in part by its perceived  
22  
23 value by clinicians, regulators, and payors, particularly in light of the limitations of trial  
24  
25 populations.<sup>19</sup> As RWE is increasingly used to refine care standards through clinical, regulatory,  
26  
27 and reimbursement pathways, its accuracy has come under increased scrutiny. This is  
28  
29 particularly important for complex medical conditions, such as HF.<sup>20</sup> Accordingly, in this  
30  
31 analysis, chart abstraction was used to quantitatively evaluate traditional and advanced  
32  
33 approaches to define HF-specific data elements. This enabled rigorous evaluation of whether  
34  
35 commonly used techniques are sufficiently accurate for observational studies, comparative  
36  
37 effectiveness research, and post-approval safety studies.  
38  
39  
40  
41  
42

43 In this study, 1) the use of an advanced, AI-based approach consistently identified HF  
44  
45 phenotypes (i.e., HF<sub>r</sub>EF, HF<sub>m</sub>rEF, and HF<sub>p</sub>EF) more accurately than a traditional approach; 2)  
46  
47 common HF symptoms and comorbid conditions were consistently and accurately identified  
48  
49 using an advanced approach; and 3) medications for HF were accurately identified using both  
50  
51 advanced and traditional approaches. While studies have previously leveraged an AI-based  
52  
53 approach to identify patients with HF,<sup>21-24</sup> the findings presented here highlight the discrepancy  
54  
55  
56  
57

1  
2  
3 between traditional EHR query methods and an AI-based approach standardized against a  
4 manual reference. Given that the accuracy of the data set and appropriateness of the applied  
5 technology are not tested in many RWE studies, there is a high potential for error.<sup>25-28</sup> The  
6 current findings highlight this while also reinforcing the impact that specific AI technologies  
7 (e.g., NLP vs. NLP plus inference) can have on phenotype generation and study validity.  
8  
9

10  
11  
12 Accurate phenotyping is paramount in any RWE study that includes HF patients. With varying  
13 etiologies and multiple phenotypes, HF is a clinically diverse syndrome, with outcomes that may  
14 vary between and even within subgroups.<sup>29,30</sup> In addition, HF patients may have different  
15 trajectories, highlighting some of the limitations of using structured data. For example, LVEF  
16 may fluctuate throughout a patient's disease course, with some patients experiencing recovery of  
17 their LVEF with the use of guideline-directed medical therapy. Accordingly, accurate  
18 phenotyping of HF patients usually requires the incorporation of data that crosses clinical  
19 encounters. In addition, although symptoms are an essential reflection of clinical status, they are  
20 poorly captured in structured data. Suboptimal recognition of comorbidities like valvular heart  
21 disease can also impact disease trajectory and risk for future cardiovascular events.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 The findings presented here represent an important advance for RWE studies that include HF  
41 patients. Notably, the only way to ascertain comparative accuracy between data sources and  
42 technologies in a domain is to test it. Accuracy consists of both recall and precision, and in the  
43 case of many health conditions, recall can fall below 50% when one relies solely upon the  
44 problem list.<sup>31,32</sup>  
45  
46  
47  
48  
49  
50

51  
52 In the current study, use of the F<sub>1</sub> score enabled analysis of both precision and recall. Despite  
53 availability of SNOMED codes for HF<sub>r</sub>EF and HF<sub>p</sub>EF, along with a similar code for HF<sub>m</sub>rEF,  
54  
55  
56  
57

1  
2  
3 such codes were rarely included. Documentation of a HF code using structured data was only  
4  
5 found 46.9% of the time when there was clear evidence of HF in the chart. The low accuracy of  
6  
7 structured data for disease subtypes may, at least partially, relate to how the data is likely to be  
8  
9 used. A physician may look within notes to understand HF subtype. Information entered into  
10  
11 problem lists and claims may be more to provide a high-level understanding of disease burden.  
12  
13 Granular billing codes may be a low priority for physicians if claims are reimbursed with the  
14  
15 non-granular HF code. Furthermore, because addition of diagnoses to the problem list is not a  
16  
17 requirement, the problem list may not be specific or updated. This contrasts with clinical notes,  
18  
19 where detailed documentation is usually performed to communicate a care plan and is a medical-  
20  
21 legal requirement.  
22  
23  
24  
25  
26

27 When low-accuracy and non-granular data are utilized, there are several potential consequences.  
28  
29 Missingness can result in selection bias, particularly if sicker patients have more frequent  
30  
31 encounters, higher rates of specialty care, and more complete documentation. Depending on the  
32  
33 study question, use of structured data alone to identify certain subgroups may be inadvisable,  
34  
35 since these data have a low recall for specific clinical concepts such as ST-elevation myocardial  
36  
37 infarction and HFpEF.<sup>33</sup> Even advanced approaches (e.g., NLP) may result in poor accuracy, as  
38  
39 illustrated in this study, where HFpEF required AI-based inference for proper identification.  
40  
41 Collectively, this highlights that not all data sources and technologies are the same; therefore,  
42  
43 accuracy testing may be required for rigorous RWE generation.<sup>34</sup> Furthermore, given the growth  
44  
45 in RWE to support new drug indications, post-marketing surveillance, and decision-making  
46  
47 regarding reimbursement, such inaccuracies may have a profound impact on large numbers of  
48  
49 patients.  
50  
51  
52  
53  
54  
55  
56  
57

1  
2  
3 Even though standard dictionaries and clinical terms related to cardiovascular medicine were  
4 used, there is a need to test the two analytic methods using different EHRs across a broader set of  
5  
6 community and referral practices. With numerous EHRs available and practitioner-to-  
7  
8 practitioner variability in documentation accuracy, efforts like the one described here represent  
9  
10 an important means of strengthening data quality.  
11  
12  
13

14  
15  
16 Importantly, this study has several limitations. First, data from a single health system was used  
17  
18 and results may not be generalizable to other populations. Second, the study protocol required  
19  
20 the selection of patients enriched with cardiovascular disease to make the study feasible, with  
21  
22 manual chart abstraction conducted to ensure the accuracy of results. While selection criteria  
23  
24 were applied to both structured and unstructured data, it is possible that this could have biased  
25  
26 results in a way that favored structured data since a larger proportion of patients with HF on the  
27  
28 problem list may have been included than if the sample had been created randomly. In addition,  
29  
30 the specific filters used likely led to a higher-than-expected proportion of HF<sub>r</sub>EF patients  
31  
32 (compared to those with HF<sub>m</sub>rEF and HF<sub>p</sub>EF). Second, the study required laborious manual  
33  
34 annotation of thousands of records. Such a sample size is adequate for high-prevalence  
35  
36 conditions, but would likely require adjustment for low-prevalence conditions with low concept  
37  
38 occurrence rates. Finally, the study did not include clinical outcome assessment; rather, it was  
39  
40 designed to compare data sources and processing methods.  
41  
42  
43  
44  
45

## 46 47 **Conclusion**

48  
49  
50 As RWE is increasingly used to analyze patient subgroups, inform clinical decision-making, and  
51  
52 influence regulatory and reimbursement decisions, data reliability and evidence validity are of  
53  
54 critical importance. Use of a traditional approach was associated with low data accuracy. While  
55  
56  
57



1  
2  
3 much greater accuracy was observed with AI-based methods, it depended upon the technology  
4 utilized. These findings highlight the importance of using data fit-for-purpose to the research  
5 question posed. In addition, they suggest that accuracy testing should be part of any EHR-based  
6 study that includes HF patients. Finally, unstructured data and a technology-based approach to  
7 data extraction may be required in some studies to achieve sufficient accuracy, depending upon  
8 the clinical assertion being tested.  
9  
10  
11  
12  
13  
14  
15

### 16 **Acknowledgments**

17  
18 We are grateful for comments from Jacob Abraham, MD, Medical Director at Providence Heart  
19 Institute's Center for Advanced Heart Disease, and Yuri Quintana, MD, Chief of, Division of  
20 Clinical Informatics at the Beth Israel Deaconess Medical Center. Editorial support was provided  
21 by Liam Gillies, Ph.D., CMPP, of Cactus Life Sciences (part of Cactus Communications),  
22 funded by Amgen Inc.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Contributors

ARG and DR drafted the manuscript. ARG, KLM, RED, DR, and TJG critically reviewed the manuscript. ARG, RED, DR, and TJG provided clinical insight.

## Funding

A research grant supported this work from Amgen Inc. DR was partly supported by the US Food and Drug Administration (FDA) under Award Number IIP-2024958 and the National Center for Advancing Translational Sciences of the NIH under Award Number R44TR002437. The content is solely the responsibility of the authors and does not necessarily represent the official views of Amgen, the FDA, or the NIH.

### **Competing interests**

KLM and RED are employees and stockholders of Amgen Inc. DR is an employee and stockholder of Verantoss, Inc. ARG has received research support from Abbott and TJG has no competing interests to declare.

### **Ethics Approval**

This study has been independently reviewed and accepted for exemption in accordance with 45 CFR 46.101(b)(4ii).

### **Provenance and peer review**

Not commissioned, externally peer reviewed.

### **Data sharing statement**

No additional data are available.

### **Supplemental Materials**

Supplemental Methods

**REFERENCES**

1. Thomas H, Diamond J, Vieco A, Chaudhuri S, Shinnar E, Cromer S, Perel P, Mensah GA, Narula J, Johnson CO, et al. Global atlas of cardiovascular disease 2000-2016: the path to prevention and control. *Glob Heart*. 2018;13:143-163. doi: 10.1016/j.gheart.2018.09.511
2. Nichols M, Townsend N, Scarborough P, Rayner M. Cardiovascular disease in Europe 2014: epidemiological update. *Eur Heart J*. 2014;35:2950-2959. doi: 10.1093/eurheartj/ehu299
3. McMurray JJ, Packer M, Desai AS, Gong J, Lefkowitz MP, Rizkala AR, Rouleau JL, Shi VC, Solomon SD, Swedberg K, et al. Angiotensin-neprilysin inhibition versus enalapril in heart failure. *N Engl J Med*. 2014;371:993-1004. doi: 10.1056/NEJMoa1409077
4. McMurray JJV, Solomon SD, Inzucchi SE, Køber L, Kosiborod MN, Martinez FA, Ponikowski P, Sabatine MS, Anand IS, Bělohávek J, et al. Dapagliflozin in patients with heart failure and reduced ejection fraction. *N Engl J Med*. 2019;381:1995-2008. doi: 10.1056/NEJMoa1911303
5. Packer M, Anker SD, Butler J, Filippatos G, Pocock SJ, Carson P, Januzzi J, Verma S, Tsutsui H, Brueckmann M, et al. Cardiovascular and renal outcomes with empagliflozin in heart failure. *N Engl J Med*. 2020;383:1413-1424. doi: 10.1056/NEJMoa2022190.
6. Swedberg K, Komajda M, Böhm M, Borer JS, Ford I, Dubost-Brama A, Lerebours G, Tavazzi L, SHIFT Investigators. Ivabradine and outcomes in chronic heart failure (SHIFT): a randomised placebo-controlled study. *Lancet*. 2010;376:875-885. doi: 10.1016/S0140-6736(10)61198-1

- 1  
2  
3 7. Stone GW, Lindenfeld J, Abraham WT, Kar S, Lim DS, Mishell JM, Whisenant B,  
4  
5 Grayburn PA, Rinaldi M, Kapadia SR, et al. Transcatheter mitral-valve repair in patients  
6  
7 with heart failure. *N Engl J Med*. 2018;379:2307-2318. doi: 10.1056/NEJMoa1806640  
8  
9
- 10 8. Shah KS, Xu H, Matsouaka RA, Bhatt DL, Heidenreich PA, Hernandez AF, Devore AD,  
11  
12 Yancy CW, Fonarow GC. Heart failure with preserved, borderline, and reduced ejection  
13  
14 fraction: 5-year outcomes. *J Am Coll Cardiol*. 2017;70:2476-2486.  
15  
16
- 17 9. H.R.34 - 21st Century Cures Act of 2016. Public Law No. 114-255. Section 3022.  
18  
19 Available at: <https://www.congress.gov/bill/114th-congress/house-bill/34/>. Accessed  
20  
21 May 4, 2020.  
22  
23
- 24 10. Pulini AA, Caetano GM, Clautiaux H, Vergeron L, Pitts PJ, Katz G. Impact of real-world  
25  
26 data on market authorization, reimbursement decision & price negotiation [published  
27  
28 online ahead of print August 28, 2020]. *Ther Innov Regul Sci*. 2020. doi:  
29  
30 10.1007/s43441-020-00208-1  
31  
32
- 33 11. Hernandez-Boussard T, Monda KL, Crespo BC, Riskin D. Real world evidence in  
34  
35 cardiovascular medicine: ensuring data validity in electronic health record-based studies.  
36  
37 *J Am Med Inform Assoc*. 2019;26:1189-1194. doi: 10.1093/jamia/ocz119.  
38  
39
- 40 12. McCormick N, Lacaille D, Bhole V, Avina-Zubieta JA. Validity of heart failure  
41  
42 diagnoses in administrative databases: a systematic review and meta-analysis. *PLoS One*.  
43  
44 2014;9:e104519. doi: 10.1371/journal.pone.0104519.  
45  
46
- 47 13. Alqaisi F, Williams LK, Peterson EL, Lanfear DE. Comparing methods for identifying  
48  
49 patients with heart failure using electronic data sources. *BMC Health Serv Res*.  
50  
51 2009;9:237. doi: 10.1186/1472-6963-9-237  
52  
53  
54  
55  
56  
57

- 1  
2  
3 14. Xu Y, Lee S, Martin E, D'souza AG, Doktorchik CTA, Jiang J, Lee S, Eastwood CA,  
4  
5 Fine N, Hemmelgarn B, et al. Enhancing ICD-code-based case definition for heart failure  
6  
7 using electronic medical record data. *J Card Fail*. 2020;26:610-617. doi:  
8  
9 10.1016/j.cardfail.2020.04.003  
10  
11  
12 15. <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world->  
13  
14 [data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world-data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory)  
15  
16  
17 16. Cunningham H, Tablan V, Roberts A, Bontcheva K. Getting more out of biomedical  
18  
19 documents with GATE's full lifecycle open source text analytics. *PLoS Comput Biol*.  
20  
21 2013;9(2):e1002854.  
22  
23  
24 17. Van Rijsbergen CJ. Information Retrieval. 2nd ed. Butterworth-Heinemann. 1979.  
25  
26 18. Bozkurt B, Coats AJ, Tsutsui H, et al. Universal Definition and Classification of Heart  
27  
28 Failure: A Report of the Heart Failure Society of America, Heart Failure Association of  
29  
30 the European Society of Cardiology, Japanese Heart Failure Society and Writing  
31  
32 Committee of the Universal Definition of Heart Failure. *J Card Fail*. 2021 Mar 1:S1071-  
33  
34 9164(21)00050-6. doi: 10.1016/j.cardfail.2021.01.022.  
35  
36  
37 19. Lim YMF, Molnar M, Vaartjes I, et al. Generalizability of randomized controlled trials in  
38  
39 heart failure with reduced ejection fraction. *Eur Heart J Qual Care Clin Outcomes*.  
40  
41 2022;8(7):761-769.  
42  
43  
44 20. Quach S, Blais C, Quan H. Administrative data have high variation in validity for  
45  
46 recording heart failure. *Can J Cardiol* 2010;26:306–12.  
47  
48  
49 21. Bielinski SJ, Pathak J, Carrell DS, et al. A Robust e-Epidemiology Tool in Phenotyping  
50  
51 Heart Failure with Differentiation for Preserved and Reduced Ejection Fraction: the  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Electronic Medical Records and Genomics (eMERGE) Network. *J Cardiovasc Transl*  
4  
5 *Res.* 2015 Nov;8(8):475-83.  
6  
7  
8 22. Blecker S, Katz SD, Horwitz LI, Kuperman G, Park H, Gold A, Sontag D. Comparison of  
9  
10 approaches for heart failure case identification from electronic health record data. *JAMA*  
11  
12 *Cardiol.* 2016;1:1014-1020. doi: 10.1001/jamacardio.2016.3236  
13  
14  
15 23. Ng K, Steinhubl SR, deFilippi C, Dey S, Stewart WF. Early detection of heart failure  
16  
17 using electronic health records: practical implications for time before diagnosis, data  
18  
19 diversity, data quantity, and data density. *Circ Cardiovasc Qual Outcomes.* 2016;9:649-  
20  
21 658. doi: 10.1161/CIRCOUTCOMES.116.002797  
22  
23  
24 24. Tison GH, Chamberlain AM, Pletcher MJ, Dunlay SM, Weston SA, Killian JM, Olgin  
25  
26 JE, Roger VL. Identifying heart failure using EMR-based algorithms. *Int J Med Inform.*  
27  
28 2018;120:1-7. doi: 10.1016/j.ijmedinf.2018.09.016  
29  
30  
31 25. Khand AU, Shaw M, Gemmel I, Cleland JG. Do discharge codes underestimate  
32  
33 hospitalisation due to heart failure? Validation study of hospital discharge coding for  
34  
35 heart failure. *Eur J Heart Fail* 2005;7: 792–797.  
36  
37  
38 26. Merry AH, Boer JM, Schouten LJ, Feskens EJ, Verschuren WM, et al. Validity of  
39  
40 coronary heart diseases and heart failure based on hospital discharge and mortality data in  
41  
42 the Netherlands using the cardiovascular registry Maastricht cohort study. *Eur J*  
43  
44 *Epidemiol* 2009;24: 237–247.  
45  
46  
47 27. U.S. Food and Drug Administration. Framework for FDA's Real-World Evidence  
48  
49 Program. 2018. Available at: <https://www.fda.gov/media/120060/download>. Accessed  
50  
51 July 26, 2020.  
52  
53  
54  
55  
56  
57

- 1  
2  
3 28. Parsons A, McCullough C, Wang J, Shih S. Validity of electronic health record-derived  
4 quality measurement for performance monitoring. *J Am Med Inform Assoc.* 2012;19:604-  
5 609. doi: 10.1136/amiajnl-2011-000557  
6  
7  
8  
9  
10 29. Kao DP, Lewsey JD, Anand IS, et al. Characterization of subgroups of heart failure  
11 patients with preserved ejection fraction with possible implications for prognosis and  
12 treatment response. *Eur J Heart Fail.* 2015;17(9):925-35.  
13  
14  
15  
16 30. Uijl A, Savarese G, Vaartjes I, et al. Identification of distinct phenotypic clusters in heart  
17 failure with preserved ejection fraction. *Eur J Heart Fail.* 2021;23(6):973-982.  
18  
19  
20  
21 31. Luna D, Franco M, Plaza C, Otero C, Wassermann S, Gambarte ML, Giunta D, de Quirós  
22 FGB. Accuracy of an electronic problem list from primary care providers and specialists.  
23 *Stud Health Technol Inform.* 2013;192:417-421.  
24  
25  
26  
27 32. Singer A, Yakubovich S, Kroeker AL, Dufault B, Duarte R, Katz A. Data quality of  
28 electronic medical records in Manitoba: do problem lists accurately reflect chronic  
29 disease billing diagnoses? *J Am Med Inform Assoc.* 2016;23:1107-1112. doi:  
30 10.1093/jamia/ocw013  
31  
32  
33  
34  
35  
36  
37 33. Makam AN, Lanham HJ, Batchelor K, Samal L, Moran B, Howell-Stampley T, Kirk L,  
38 Cherukuri M, Santini N, Leykum LK, et al. Use and satisfaction with key functions of a  
39 common commercial electronic health record: a survey of primary care providers. *BMC*  
40 *Med Inform Decis Mak.* 2013;13:86. doi: 10.1186/1472-6947-13-86.  
41  
42  
43  
44  
45  
46  
47 34. Ingelsson E, Arnlov J, Sundstrom J, Lind L. The validity of a diagnosis of heart failure in  
48 a hospital discharge register. *Eur J Heart Fail* 2005;7:787-791.  
49  
50  
51  
52  
53  
54  
55  
56  
57



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

Table 1. Prespecified heart failure-specific concepts extracted from the electronic health record.

High Priority Conditions	Comorbidities	Symptoms	Findings	Medications
Congestive HF	Myocardial infarction	Angina	LVEF	Carvedilol
HF with reduced EF	Atrial fibrillation	Chest pain		Lisinopril
HF with mid-range EF	Aortic regurgitation	Dyspnea		Metoprolol
HF with preserved EF	Mitral regurgitation	Fatigue		Furosemide
	Tricuspid regurgitation	Palpitations		

HF, heart failure; EF, ejection fraction; LVEF, left ventricular ejection fraction.

Peer review only

### EHR data source and processing

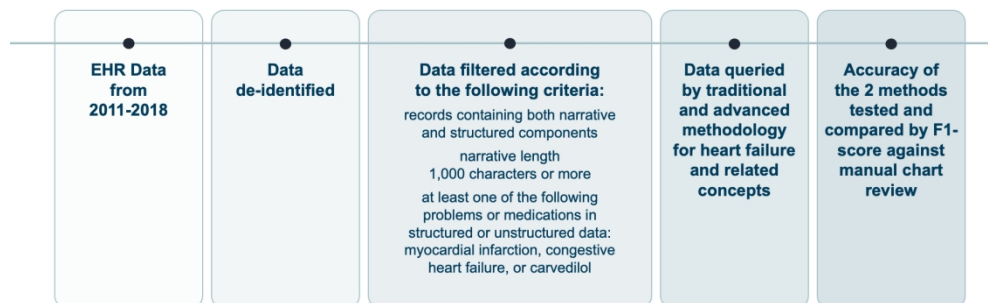


Figure 1: Electronic Health Record data source and processing.

338x190mm (144 x 144 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

### EHR data source and processing

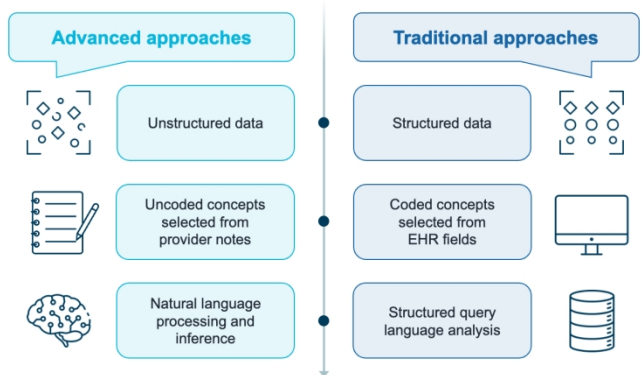


Figure 2: Comparison of traditional and advanced real-world evidence approaches. EHR, electronic health record.

338x190mm (144 x 144 DPI)

### F<sub>1</sub> score by subcategories of heart failure

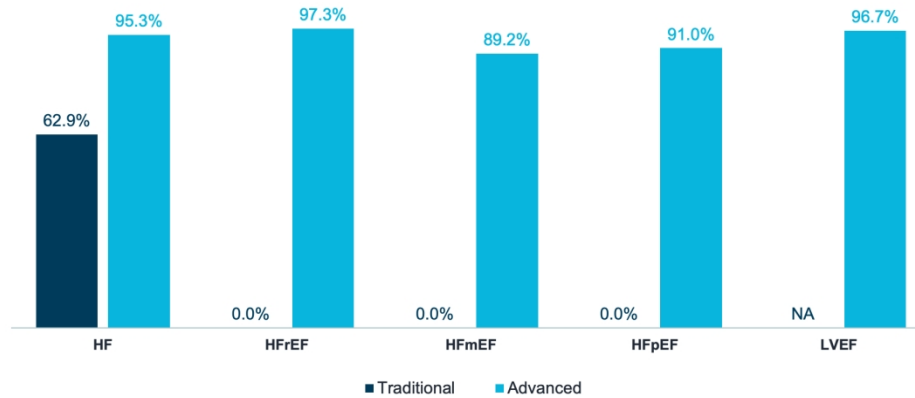
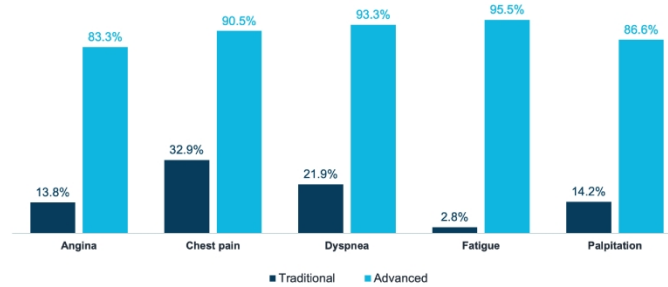
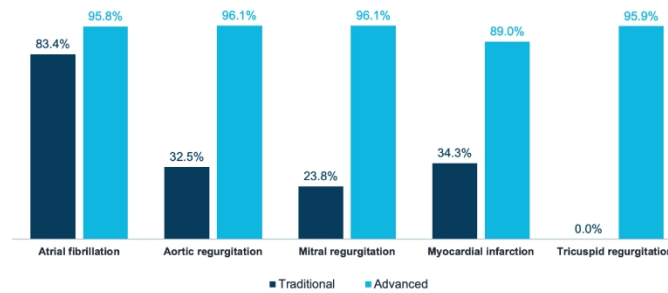
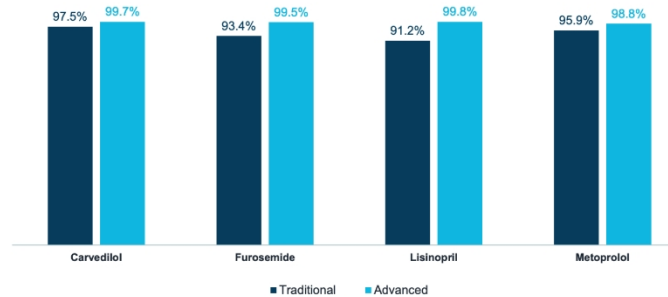


Figure 3: F1 scores for heart failure diagnoses. \*F1-score could not be calculated due to lack of data for precision. †Structured data recall is not applicable for ejection fraction because no code was available within the problem list. HF, heart failure; HFmEF, heart failure with mildly-reduced ejection fraction; HFpEF, heart failure with preserved ejection fraction; HFReEF, heart failure with reduced ejection fraction; LVEF, left ventricular ejection fraction; 0% reflects a measured value and indicates the availability of the diagnosis code in the EHR dropdown versus N/A, not applicable, which refers to a diagnosis without available code in the relevant codeset.

338x190mm (144 x 144 DPI)

**F<sub>1</sub> score by symptoms****F<sub>1</sub> score by type of comorbidity****F<sub>1</sub> score by medication type**

F1 scores for (A) symptoms, (B) medications, and (C) comorbid conditions. \*F1 score could not be calculated due to a lack of data for precision. N/A, not applicable.

548x904mm (118 x 118 DPI)

Supplementary Table 1. Cohort identification of heart failure diagnoses, left ventricular ejection fraction, heart failure medications, symptoms, and comorbid cardiovascular conditions

	Traditional approach			Advanced approach			Concept occurrence	Encounter occurrence	P-value
	Recall, %	Precision, %	F <sub>1</sub> score, %	Recall, %	Precision, %	F <sub>1</sub> score, %			
<b>HF diagnosis</b>									
HF	46.9	95.4	62.9	96.0	94.7	95.3	265	155	<0.001
HFrEF	0	N/A*	N/A <sup>†</sup>	94.8	100.0	97.3	382	124	N/A <sup>§</sup>
HFmrEF	0	N/A*	N/A <sup>†</sup>	80.4	100.0	89.2	62	35	N/A <sup>§</sup>
HFpEF	0	N/A*	N/A <sup>†</sup>	83.5	100.0	91.0	103	55	N/A <sup>§</sup>
<b>LVEF</b>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	93.7	100.0	96.7	677	238	N/A <sup>§</sup>
<b>HF medications</b>									
Carvedilol	95.1	100.0	97.5	99.7	99.7	99.7	407	141	<0.001
Furosemide	87.7	100.0	93.4	99.3	99.8	99.5	1572	371	0.116
Lisinopril	83.9	100.0	91.2	99.7	99.9	99.8	1068	386	<0.001
Metoprolol	92.2	100.0	95.9	97.7	100.0	98.8	1370	397	<0.001
<b>Symptoms</b>									

Angina	7.8	60.0	13.8	84.4	82.3	83.3	265	155	<0.00 1
Chest pain	21.4	70.8	32.9	95.4	86.1	90.5	2332	756	<0.00 1
Dyspnea	12.7	78.2	21.9	94.7	92.0	93.3	4474	832	<0.00 1
Fatigue	1.4	75.0	2.8	96.5	94.5	95.5	1711	371	<0.00 1
Palpitation	8.2	52.9	14.2	90.9	82.6	86.6	896	493	<0.00 1
<b>Comorbid cardiovascular conditions</b>									
Atrial fibrillation	72.2	98.7	83.4	93.0	98.7	95.8	1214	222	<0.00 1
Aortic regurgitation	19.4	100.0	32.5	92.5	100.0	96.1	153	90	<0.00 1
Mitral regurgitation	13.5	97.1	23.8	92.8	99.6	96.1	483	185	<0.00 1
Myocardial infarction	21.1	90.9	34.3	95.5	83.4	89.0	1220	578	<0.00 1
Tricuspid regurgitation	0	N/A*	N/A†	92.2	100.0	95.9	162	78	N/A§

\*These elements did not occur when using the traditional approach. †F<sub>1</sub> scores could not be calculated due to a lack of data for precision. ‡Structured data recall is not applicable for ejection fraction because there was no code available within the problem list. §P-value could not be calculated due to the unavailability of F<sub>1</sub> scores for the traditional approach. P-values are derived from the chi-square test.

1  
2  
3 HF, heart failure; HFmrEF, heart failure with mildly reduced ejection fraction; HFpEF, heart failure with  
4 preserved ejection fraction; HFrEF, heart failure with reduced ejection fraction; LVEF, left ventricular  
5 ejection fraction; N/A, not applicable.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## SUPPLEMENTAL MATERIAL

### Reference Standard

Traditional and advanced approaches were tested against a reference standard for physician encounters. The reference standard consisted of an independent review, with manual annotation of relevant HF-specific features, including 19 unique HF-specific concepts. For each encounter, two independent clinical annotators labeled each concept and all metadata for that concept. For example, an annotator might mark the text "DOE over last month" as dyspnea on exertion, experienced = true, current = true, relative date = 1 month. Concept occurrence was defined as the sum of all concept occurrences, allowing for multiple occurrences per encounter. Encounter occurrence was defined as the number of encounters with at least one occurrence of the concept.

Given that many concepts, such as LVEF are specific to a point in time, concepts were tested at the encounter level. For example, if a patient had an LVEF of 30% in an encounter, the data extraction would only be annotated as correct if it identified "LVEF 30%" in that specific encounter. This reference standard was used to determine accuracy of automated extracted data and structured data. Specifically, this reference standard was used to calculate recall and precision for these individual features for traditional and advanced approaches.

# BMJ Open

## A Retrospective Comparison of Traditional and Artificial-Intelligence Based Heart Failure Phenotyping in a US Health System to Enable Real-World Evidence

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2023-073178.R2
Article Type:	Original research
Date Submitted by the Author:	12-Jul-2023
Complete List of Authors:	Garan, Arthur Reshad; Harvard Medical School Monda , Keri; Amgen Inc Dent-Acosta, Ricardo ; Amgen Inc Riskin , Daniel; Verantos Gluckman, Ty; Providence St Joseph Health
<b>Primary Subject Heading</b>:	Health informatics
Secondary Subject Heading:	Cardiovascular medicine
Keywords:	Heart failure < CARDIOLOGY, Cardiac Epidemiology < CARDIOLOGY, Health informatics < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

1  
2  
3 **A Retrospective Comparison of Traditional and Artificial-Intelligence Based Heart Failure**  
4  
5 **Phenotyping in a US Health System to Enable Real-World Evidence**  
6  
7

8  
9 **Short title:** Real-World Evidence in Heart Failure  
10

11  
12 **Authors:** A. Reshad Garan, MD, MS; Keri L. Monda, Ph.D.; Ricardo E. Dent-Acosta; MD;  
13  
14 Dan Riskin, MD; Ty J. Gluckman, MD, MHA  
15  
16

17  
18 **Affiliations:** Department of Medicine, Division of Cardiology, Beth Israel Deaconess Medical  
19  
20 Center, Harvard Medical School, Boston, MA, USA (A.R.G.). The Center for Observational  
21  
22 Research and Medical Affairs, Amgen Inc., Thousand Oaks, CA, USA (K.L.M., R.E.D.-A.).  
23  
24 Department of Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, NC,  
25  
26 USA (K.L.M.). Verantos, Inc., Menlo Park, CA, USA (D.R.). Department of Surgery, Stanford  
27  
28 University School of Medicine, Stanford, CA, USA (D.R.). Center for Cardiovascular Analytics,  
29  
30 Research and Data Science (CARDS), Providence Heart Institute, Providence St. Joseph Health,  
31  
32 Portland, OR, USA (T.J.G.)  
33  
34  
35  
36

37  
38 **Address for correspondence:** A. Reshad Garan, MD, MS, Department of Medicine, Division of  
39  
40 Cardiology, Beth Israel Deaconess Medical Center, 185 Pilgrim Road Deaconess 3, Boston, MA  
41  
42 02215, USA. Phone: 617/632-7737. Fax: 617/632-7620. Email: [agaran@bidmc.harvard.edu](mailto:agaran@bidmc.harvard.edu)  
43  
44  
45

46 **Word count:** 3136 words  
47  
48

49 **Keywords:** artificial intelligence; heart failure; phenotype; real-world evidence; electronic health  
50  
51 record.  
52  
53

54  
55 **Abstract**  
56  
57

1  
2  
3 **Objective:** Quantitatively evaluate the quality of data underlying real-world evidence (RWE) in  
4 heart failure (HF).  
5  
6

7 **Design:** Retrospective comparison of accuracy in identifying HF patients and phenotypic  
8 information was made using traditional (i.e., structured query language applied to structured  
9 EHR data) and advanced (i.e., AI applied to unstructured EHR data) RWE approaches. The  
10 performance of each approach was measured by the harmonic mean of precision and recall (F1  
11 score) using manual annotation of medical records as a reference standard.  
12  
13  
14  
15  
16  
17  
18

19 **Setting:** EHR data from a large academic healthcare system in North America between 2015 and  
20 2019, with an expected catchment of approximately 500,000 patients.  
21  
22  
23

24 **Population:** 4288 encounters for 1155 patients aged 18 to 85 years, with 472 patients identified  
25 as having HF.  
26  
27

28 **Outcome measures:** HF and associated concepts, such as comorbidities, left ventricular ejection  
29 fraction, and selected medications.  
30  
31  
32  
33

34 **Results:** The average F1 scores across 19 HF-specific concepts were 49.0% and 94.1% for the  
35 traditional and advanced approaches, respectively ( $P < 0.001$  for all concepts with available data).  
36 The absolute difference in F<sub>1</sub> score between approaches was 45.1% (98.1% relative increase in  
37 F<sub>1</sub> score using the advanced approach). The advanced approach achieved superior F1 scores for  
38 HF presence, phenotype, and associated comorbidities. Some phenotypes, such as HFpEF,  
39 revealed dramatic differences in extraction accuracy based on technology applied, with a 4.9%  
40 F<sub>1</sub> score when using natural language processing (NLP) alone and a 91.0% F<sub>1</sub> score when using  
41 NLP plus AI-based inference.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 **Conclusions:** A traditional RWE generation approach resulted in low data quality in HF patients.  
4  
5 While an advanced approach demonstrated high accuracy, the results varied dramatically based  
6  
7 on extraction techniques. For future studies, advanced approaches and accuracy measurement  
8  
9 may be required to ensure data are fit-for-purpose.  
10  
11  
12

### 13 **Strengths and limitations of this study**

- 14  
15  
16  
17 • Using RWE for HF patients requires demonstrating that the data source and technologies  
18  
19 result in accurate data.
- 20  
21  
22 • Natural language processing alone lacked context from the longitudinal record, limiting  
23  
24 phenotype identification and study validity.
- 25  
26  
27 • Findings suggest that advanced methods can enable high-validity RWE for heart failure  
28  
29 patients.
- 30  
31 • The use of data from a single healthcare system may limit generalizability to other  
32  
33 populations.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

## INTRODUCTION

Heart failure (HF) is a major public health problem with significant associated morbidity, mortality, and cost.<sup>1,2</sup> Despite the availability of novel drugs and devices, morbidity and mortality in HF rivals many malignancies, with a 5-year survival rate as low as 50%.<sup>3-8</sup>

Randomized controlled trials (RCTs) have traditionally been used to assess the safety and efficacy of new therapies and represent a cornerstone for regulatory approval. However, RCTs are frequently conducted in highly selected populations, typically younger, healthier, and less diverse than patients treated in clinical practice. Furthermore, such trials often include patients with an established HF diagnosis, receiving guideline-directed medical therapy at tertiary centers, and may not represent the broader HF population. Because HF is a clinically heterogeneous syndrome with numerous etiologies and phenotypes, studying this population can be particularly difficult.

Real-world evidence (RWE) has held promise as a potential means to assess therapeutic benefit outside of clinical trials, with sufficient power to characterize therapeutic impact in HF subgroups. Accordingly, RWE can complement RCTs, extending the findings to patient populations that may have been excluded from or insufficiently enrolled in pivotal trials. To accelerate these and similar precision medicine goals, the 21st Century Cures Act was passed in 2016, which required the United States Food and Drug Administration (FDA) to develop guidance supporting the use of RWE in new drug indications and post-marketing surveillance.<sup>9</sup> In addition, payors have increasingly utilized RWE to inform reimbursement decisions and are increasingly demanding credible evidence.<sup>10</sup>

1  
2  
3 Not surprisingly, the quality of RWE hinges on how well real-world data are collected,  
4 processed<sup>11</sup>, and used to inform study questions. Such is the case in HF, where accurate  
5 identification of patients in administrative and other structured data sets is an ongoing focus.<sup>12-14</sup>  
6  
7 Traditional methods of identifying HF patients rely on querying diagnosis codes and structured  
8 data in the electronic health record (EHR) or medical claims. Conversely, artificial intelligence  
9 (AI) applied to unstructured data represents a novel method of analyzing the medical record.  
10  
11 Because of the importance of data reliability in RWE and the potential to use unstructured data to  
12 achieve data enrichment<sup>15</sup>, we sought to compare the accuracy achieved by traditional RWE  
13 methods versus advanced AI approaches in identifying a range of HF-specific data elements  
14 from the medical record.  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

## 27 **METHODS**

28  
29  
30 The study design is outlined in Figure 1. Varied data sources and applied technologies were used  
31 to assess data reliability in patients with risk factors for HF. Leveraging manual chart abstraction  
32 as the reference standard, comparisons were made between the two methods. The first method  
33 used structured EHR data (e.g., diagnosis codes and problem lists) and standard query  
34 techniques, defined as the 'traditional approach'. The second used unstructured EHR data (e.g.,  
35 narratives from primary care and specialty notes) and AI techniques, described as the 'advanced  
36 approach' (Figure 1). The primary objective was measurement of the accuracy of identified HF-  
37 specific elements using traditional and advanced approaches. We hypothesized that the advanced  
38 approach would better identify key HF-specific elements than the traditional approach. Data  
39 were deidentified before study initiation, and the study was determined not to be human subjects  
40 research. Both natural language processing (NLP) and machine-learned inference technologies  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57



1  
2  
3 used in the advanced approach were provided by Verantos, Inc. (Menlo Park, CA, USA). The  
4  
5 core of AI is a deterministic NLP layer. This layer is built on top of the GATE NLP  
6  
7 architecture.<sup>16</sup> The architecture is used to construct a flexible pipeline for processing incoming  
8  
9 text against English language syntactical rules augmented with a lexicon based on a clinical  
10  
11 vocabulary. The AI-based inference was applied during data processing. Millions of machine-  
12  
13 learned and manually curated associations enable disambiguation and identification of clinically  
14  
15 relevant concepts. As an example of AI-based inference, a patient with HF on the problem list  
16  
17 and a narrative encounter describing “EF 60%” would not be interpreted by NLP as having HF  
18  
19 with preserved ejection fraction (HFpEF) since the text does not have sufficient information to  
20  
21 identify this condition. On the other hand, AI-based inference would infer HFpEF based on  
22  
23 disparate information in the record.  
24  
25  
26  
27  
28

### 29 **EHR Data Source and Processing**

30  
31  
32  
33 EHR data from primary care encounters between 2011 and 2018 were deidentified and securely  
34  
35 transferred to a cloud-based server for analysis. The data set consisted of both structured data  
36  
37 (e.g., medical conditions, procedures performed, medications, and problem lists) and  
38  
39 unstructured data (e.g., narrative notes from primary care providers and specialists, telephone  
40  
41 visits, and other narrative text) (Figure 2).  
42  
43  
44

45  
46 As the study aimed to test the accuracy of different RWE approaches and not treatment  
47  
48 effectiveness, the cohort was enriched for patients with suspected HF based on comorbidities and  
49  
50 medications. Specifically, the following filters were applied: records containing both narrative  
51  
52 and structured components; narrative length 1,000 characters or more; and at least one of the  
53  
54  
55  
56  
57

1  
2  
3 following problems or medications in structured or unstructured data: myocardial infarction,  
4  
5 congestive heart failure, or carvedilol (Figure 1).  
6  
7

8  
9 A prespecified set of clinical concepts pertinent to patients with HF was extracted using  
10  
11 traditional and advanced techniques (Table 1). Problem lists were mapped to Systematized  
12  
13 Nomenclature of Medicine (SNOMED) ontology, and unadjudicated claims were mapped to  
14  
15 ICD-10 codes. Standard sets of individual codes were used to represent each concept. With the  
16  
17 advanced approach, inference incorporating pattern recognition was utilized to identify  
18  
19 potentially missing or ignored concepts within the text (e.g., HF being likely in patients with  
20  
21 dyspnea and pitting edema on a diuretic). Specifically, no narrative coding took place before the  
22  
23 AI algorithm was used; instead, it was applied directly to the narrative text and then mapped by  
24  
25 the algorithm to the SNOMED ontology. Next, manual chart abstraction using the same  
26  
27 SNOMED code set was used as a reference to assess the accuracy of the coding by the AI  
28  
29 algorithm. Engineers were blinded to validation data and its corresponding chart abstraction.  
30  
31  
32  
33  
34

### 35 **Study End Points and Statistical Analysis**

36  
37  
38 The primary endpoint was the  $F_1$  score for traditional and advanced approaches. The  $F_1$  score is  
39  
40 an accuracy measure that combines recall and precision; more specifically, it is the weighted  
41  
42 harmonic mean of these two measures. Secondary endpoints were recall (i.e., the proportion of  
43  
44 patients correctly identified as having the condition, akin to sensitivity) and precision (i.e., the  
45  
46 proportion of patients with HF and its subtypes correctly identified divided by the total number  
47  
48 of patients identified in each cohort akin to positive predictive value)<sup>17,18</sup> for the traditional and  
49  
50 advanced approaches. The reference standard used to evaluate accuracy of the traditional and  
51  
52 advanced approaches was manual chart abstraction. For each encounter, two independent clinical  
53  
54  
55  
56  
57

1  
2  
3 annotators labeled each concept and all metadata for that concept. Annotators were blinded to  
4  
5 each other's annotations, and inter-rater agreement was measured by Cohen's kappa score.  
6  
7  
8 Further description of the reference standard methodology is provided in the Supplemental  
9  
10 Material. Results were summarized using descriptive statistics, and percentages were calculated  
11  
12 for categorical variables. Differences in  $F_1$  scores between traditional and advanced approaches  
13  
14 were analyzed using the chi-square test; associated  $P$ -values were reported.  
15  
16  
17

### 18 **Patient and Public Involvement**

19  
20  
21 Data were deidentified before study initiation, and the study was determined not to be human  
22  
23 subjects research. As a result, no patients were recruited for study participation. The research  
24  
25 question and study goal of highlighting methods for improving RWE use were driven by  
26  
27 recognition that improvements in use of RWE to inform new drug indications, post-marketing  
28  
29 surveillance, and reimbursement decisions would ultimately result in patient benefit.  
30  
31  
32  
33

### 34 **RESULTS**

35  
36  
37 A total of 4288 encounters for 1155 patients were examined, of which 472 patients with HF were  
38  
39 identified. Of these, 382 had HF with reduced ejection fraction (HF<sub>r</sub>EF), 35 had HF with mildly  
40  
41 reduced ejection fraction (HF<sub>mr</sub>EF), and 55 had HF with preserved ejection fraction (HF<sub>p</sub>EF).  
42  
43  
44 The reference standard Cohen's kappa score was 0.95, suggesting high validity.  
45  
46  
47

48  
49 Supplementary Table 1 reports the  $F_1$  score, recall, and precision results achieved with both  
50  
51 approaches. Figure 3 graphically presents  $F_1$  scores for HF diagnoses and Figure 4 includes  $F_1$   
52  
53 scores for symptoms, medications, and comorbid conditions. Overall, accuracy was significantly  
54  
55 greater for the advanced approach (AI applied to unstructured EHR data) than for the traditional  
56  
57

1  
2  
3 approach (structured query language applied to structured EHR data) (Supplementary Table 1;  
4  
5 Figure 3; Figure 4), with an absolute difference of 45.1%.

6  
7  
8  
9 With the traditional approach, recall for any HF diagnosis was 46.9% (i.e., 53.1% of patients  
10  
11 with HF were missed entirely) and precision was 95.4%, resulting in an  $F_1$  score of 62.9%  
12  
13 ( $P<0.001$ ). In contrast, with the advanced approach, recall for any HF diagnosis was 96.0% and  
14  
15 precision was 94.7%, resulting in an  $F_1$ -score of 95.3% ( $P<0.001$  when  $F_1$  scores for the two  
16  
17 approaches were compared) (Supplementary Table 1; Figure 3). Among HF phenotypes, recall  
18  
19 with the advanced approach was highest with HFrEF, followed by HFpEF and HFmrEF;  
20  
21 precision was 100% for all phenotypes. With the traditional approach,  $F_1$  scores could not be  
22  
23 calculated for HFrEF, HFmrEF, and HFpEF because only less granular HF codes were used  
24  
25 (Supplementary Table 1).

26  
27  
28  
29  
30  
31 Accuracy in identifying left ventricular ejection fraction (LVEF) was similarly high with the  
32  
33 advanced approach, with an  $F_1$  score of 96.7%. Data could not be extracted for LVEF with the  
34  
35 traditional approach because no such codes were available within the EHR, nor did a mechanism  
36  
37 to encode LVEF within the problem list or unadjudicated claims exist (Supplementary Table 1;  
38  
39 Figure 3).

40  
41  
42  
43 Accurate identification of HF symptoms was greater with the advanced approach ( $P<0.001$ )  
44  
45 (Supplementary Table 1; Figure 4A). Whereas identification of commonly prescribed HF  
46  
47 medications was high with both approaches (Supplementary Table 1; Figure 4B), identification  
48  
49 of cardiovascular comorbidities was higher in all cases with the advanced approach ( $P<0.001$ )  
50  
51 (Supplementary Table 1; Figure 4C).

1  
2  
3 Data concept extraction with the advanced approach greatly depended upon the technology used.  
4  
5 For example, NLP, which ends at the sentence boundary, was only able to identify HFpEF with  
6  
7 an F<sub>1</sub> score of 4.9% because "HFpEF" or "heart failure with preserved ejection fraction" was  
8  
9 rarely written. Conversely, inference, which can find related items from the longitudinal record,  
10  
11 was able to identify both "HF" and "normal ejection fraction" as separate annotations for HFpEF  
12  
13 with an F<sub>1</sub>-score of 91.0% (Supplementary Table 1; Figure 3).  
14  
15  
16  
17

## 18 **DISCUSSION**

19  
20  
21 The utilization of RWE has grown substantially in recent years, driven in part by its perceived  
22  
23 value by clinicians, regulators, and payors, particularly in light of the limitations of trial  
24  
25 populations.<sup>19</sup> As RWE is increasingly used to refine care standards through clinical, regulatory,  
26  
27 and reimbursement pathways, its accuracy has come under increased scrutiny. This is  
28  
29 particularly important for complex medical conditions, such as HF.<sup>20</sup> Accordingly, in this  
30  
31 analysis, chart abstraction was used to quantitatively evaluate traditional and advanced  
32  
33 approaches to define HF-specific data elements. This enabled rigorous evaluation of whether  
34  
35 commonly used techniques are sufficiently accurate for observational studies, comparative  
36  
37 effectiveness research, and post-approval safety studies.  
38  
39  
40  
41  
42

43 In this study, 1) the use of an advanced, AI-based approach consistently identified HF  
44  
45 phenotypes (i.e., HF<sub>r</sub>EF, HF<sub>m</sub>rEF, and HF<sub>p</sub>EF) more accurately than a traditional approach; 2)  
46  
47 common HF symptoms and comorbid conditions were consistently and accurately identified  
48  
49 using an advanced approach; and 3) medications for HF were accurately identified using both  
50  
51 advanced and traditional approaches. While studies have previously leveraged an AI-based  
52  
53 approach to identify patients with HF,<sup>21-24</sup> the findings presented here highlight the discrepancy  
54  
55  
56  
57

1  
2  
3 between traditional EHR query methods and an AI-based approach standardized against a  
4 manual reference. Given that the accuracy of the data set and appropriateness of the applied  
5 technology are not tested in many RWE studies, there is a high potential for error.<sup>25-28</sup> The  
6 current findings highlight this while also reinforcing the impact that specific AI technologies  
7 (e.g., NLP vs. NLP plus inference) can have on phenotype generation and study validity.  
8  
9

10  
11  
12 Accurate phenotyping is paramount in any RWE study that includes HF patients. With varying  
13 etiologies and multiple phenotypes, HF is a clinically diverse syndrome, with outcomes that may  
14 vary between and even within subgroups.<sup>29,30</sup> In addition, HF patients may have different  
15 trajectories, highlighting some of the limitations of using structured data. For example, LVEF  
16 may fluctuate throughout a patient's disease course, with some patients experiencing recovery of  
17 their LVEF with the use of guideline-directed medical therapy. Accordingly, accurate  
18 phenotyping of HF patients usually requires the incorporation of data that crosses clinical  
19 encounters. In addition, although symptoms are an essential reflection of clinical status, they are  
20 poorly captured in structured data. Suboptimal recognition of comorbidities like valvular heart  
21 disease can also impact disease trajectory and risk for future cardiovascular events.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 The findings presented here represent an important advance for RWE studies that include HF  
41 patients. Notably, the only way to ascertain comparative accuracy between data sources and  
42 technologies in a domain is to test it. Accuracy consists of both recall and precision, and in the  
43 case of many health conditions, recall can fall below 50% when one relies solely upon the  
44 problem list.<sup>31,32</sup>  
45  
46  
47  
48  
49  
50

51  
52 In the current study, use of the  $F_1$  score enabled analysis of both precision and recall. Despite  
53 availability of SNOMED codes for HF<sub>r</sub>EF and HF<sub>p</sub>EF, along with a similar code for HF<sub>m</sub>rEF,  
54  
55  
56  
57

1  
2  
3 such codes were rarely included. Documentation of a HF code using structured data was only  
4  
5 found 46.9% of the time when there was clear evidence of HF in the chart. The low accuracy of  
6  
7 structured data for disease subtypes may, at least partially, relate to how the data is likely to be  
8  
9 used. A physician may look within notes to understand HF subtype. Information entered into  
10  
11 problem lists and claims may be more to provide a high-level understanding of disease burden.  
12  
13 Granular billing codes may be a low priority for physicians if claims are reimbursed with the  
14  
15 non-granular HF code. Furthermore, because addition of diagnoses to the problem list is not a  
16  
17 requirement, the problem list may not be specific or updated. This contrasts with clinical notes,  
18  
19 where detailed documentation is usually performed to communicate a care plan and is a medical-  
20  
21 legal requirement.  
22  
23  
24  
25  
26

27 When low-accuracy and non-granular data are utilized, there are several potential consequences.  
28  
29 Missingness can result in selection bias, particularly if sicker patients have more frequent  
30  
31 encounters, higher rates of specialty care, and more complete documentation. Depending on the  
32  
33 study question, use of structured data alone to identify certain subgroups may be inadvisable,  
34  
35 since these data have a low recall for specific clinical concepts such as ST-elevation myocardial  
36  
37 infarction and HF<sub>r</sub>EF.<sup>33</sup> Even advanced approaches (e.g., NLP) may result in poor accuracy, as  
38  
39 illustrated in this study, where HF<sub>p</sub>EF required AI-based inference for proper identification.  
40  
41 Collectively, this highlights that not all data sources and technologies are the same; therefore,  
42  
43 accuracy testing may be required for rigorous RWE generation.<sup>34</sup> Furthermore, given the growth  
44  
45 in RWE to support new drug indications, post-marketing surveillance, and decision-making  
46  
47 regarding reimbursement, it is imperative for clinicians to understand that such inaccuracies may  
48  
49 have a profound impact on large numbers of patients.  
50  
51  
52  
53  
54  
55  
56  
57

1  
2  
3 Even though standard dictionaries and clinical terms related to cardiovascular medicine were  
4 used, there is a need to test the two analytic methods using different EHRs across a broader set of  
5  
6 community and referral practices. With numerous EHRs available and practitioner-to-  
7  
8 practitioner variability in documentation accuracy, efforts like the one described here represent  
9  
10 an important means of strengthening data quality.  
11  
12  
13

14  
15  
16 Importantly, this study has several limitations. First, data from a single health system was used  
17  
18 and results may not be generalizable to other populations. Second, the study protocol required  
19  
20 the selection of patients enriched with cardiovascular disease to make the study feasible, with  
21  
22 manual chart abstraction conducted to ensure the accuracy of results. While selection criteria  
23  
24 were applied to both structured and unstructured data, it is possible that this could have biased  
25  
26 results in a way that favored structured data since a larger proportion of patients with HF on the  
27  
28 problem list may have been included than if the sample had been created randomly. In addition,  
29  
30 the specific filters used likely led to a higher-than-expected proportion of HF<sub>r</sub>EF patients  
31  
32 (compared to those with HF<sub>m</sub>rEF and HF<sub>p</sub>EF). Second, the study required laborious manual  
33  
34 annotation of thousands of records. Such a sample size is adequate for high-prevalence  
35  
36 conditions, but would likely require adjustment for low-prevalence conditions with low concept  
37  
38 occurrence rates. Finally, the study did not include clinical outcome assessment; rather, it was  
39  
40 designed to compare data sources and processing methods.  
41  
42  
43  
44  
45

## 46 47 **Conclusion**

48  
49  
50 As RWE is increasingly used to analyze patient subgroups, inform clinical decision-making, and  
51  
52 influence regulatory and reimbursement decisions, data reliability and evidence validity are of  
53  
54 critical importance. Use of a traditional approach was associated with low data accuracy. While  
55  
56  
57



1  
2  
3 much greater accuracy was observed with AI-based methods, it depended upon the technology  
4 utilized. These findings highlight the importance of using data fit-for-purpose to the research  
5 question posed. In addition, they suggest that accuracy testing should be part of any EHR-based  
6 study that includes HF patients. Finally, unstructured data and a technology-based approach to  
7 data extraction may be required in some studies to achieve sufficient accuracy, depending upon  
8 the clinical assertion being tested.  
9  
10  
11  
12  
13  
14  
15

### 16 **Acknowledgments**

17  
18 We are grateful for comments from Jacob Abraham, MD, Medical Director at Providence Heart  
19 Institute's Center for Advanced Heart Disease, and Yuri Quintana, MD, Chief of, Division of  
20 Clinical Informatics at the Beth Israel Deaconess Medical Center. Editorial support was provided  
21 by Liam Gillies, Ph.D., CMPP, of Cactus Life Sciences (part of Cactus Communications),  
22 funded by Amgen Inc.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Contributors

ARG and DR drafted the manuscript. ARG, KLM, RED, DR, and TJG critically reviewed the manuscript. ARG, RED, DR, and TJG provided clinical insight.

## Funding

A research grant supported this work from Amgen Inc. DR was partly supported by the US Food and Drug Administration (FDA) under Award Number IIP-2024958 and the National Center for Advancing Translational Sciences of the NIH under Award Number R44TR002437. The content is solely the responsibility of the authors and does not necessarily represent the official views of Amgen, the FDA, or the NIH.

### **Competing interests**

KLM and RED are employees and stockholders of Amgen Inc. DR is an employee and stockholder of Verantoss, Inc. ARG has received research support from Abbott and TJG has no competing interests to declare.

### **Ethics Approval**

This study has been independently reviewed and accepted for exemption in accordance with 45 CFR 46.101(b)(4ii).

### **Provenance and peer review**

Not commissioned, externally peer reviewed.

### **Data sharing statement**

No additional data are available.

### **Supplemental Materials**

Supplemental Methods

**REFERENCES**

1. Thomas H, Diamond J, Vieco A, Chaudhuri S, Shinnar E, Cromer S, Perel P, Mensah GA, Narula J, Johnson CO, et al. Global atlas of cardiovascular disease 2000-2016: the path to prevention and control. *Glob Heart*. 2018;13:143-163. doi: 10.1016/j.gheart.2018.09.511
2. Nichols M, Townsend N, Scarborough P, Rayner M. Cardiovascular disease in Europe 2014: epidemiological update. *Eur Heart J*. 2014;35:2950-2959. doi: 10.1093/eurheartj/ehu299
3. McMurray JJ, Packer M, Desai AS, Gong J, Lefkowitz MP, Rizkala AR, Rouleau JL, Shi VC, Solomon SD, Swedberg K, et al. Angiotensin-neprilysin inhibition versus enalapril in heart failure. *N Engl J Med*. 2014;371:993-1004. doi: 10.1056/NEJMoa1409077
4. McMurray JJV, Solomon SD, Inzucchi SE, Køber L, Kosiborod MN, Martinez FA, Ponikowski P, Sabatine MS, Anand IS, Bělohávek J, et al. Dapagliflozin in patients with heart failure and reduced ejection fraction. *N Engl J Med*. 2019;381:1995-2008. doi: 10.1056/NEJMoa1911303
5. Packer M, Anker SD, Butler J, Filippatos G, Pocock SJ, Carson P, Januzzi J, Verma S, Tsutsui H, Brueckmann M, et al. Cardiovascular and renal outcomes with empagliflozin in heart failure. *N Engl J Med*. 2020;383:1413-1424. doi: 10.1056/NEJMoa2022190.
6. Swedberg K, Komajda M, Böhm M, Borer JS, Ford I, Dubost-Brama A, Lerebours G, Tavazzi L, SHIFT Investigators. Ivabradine and outcomes in chronic heart failure (SHIFT): a randomised placebo-controlled study. *Lancet*. 2010;376:875-885. doi: 10.1016/S0140-6736(10)61198-1

- 1  
2  
3 7. Stone GW, Lindenfeld J, Abraham WT, Kar S, Lim DS, Mishell JM, Whisenant B,  
4  
5 Grayburn PA, Rinaldi M, Kapadia SR, et al. Transcatheter mitral-valve repair in patients  
6  
7 with heart failure. *N Engl J Med*. 2018;379:2307-2318. doi: 10.1056/NEJMoa1806640  
8  
9
- 10 8. Shah KS, Xu H, Matsouaka RA, Bhatt DL, Heidenreich PA, Hernandez AF, Devore AD,  
11  
12 Yancy CW, Fonarow GC. Heart failure with preserved, borderline, and reduced ejection  
13  
14 fraction: 5-year outcomes. *J Am Coll Cardiol*. 2017;70:2476-2486.  
15  
16
- 17 9. H.R.34 - 21st Century Cures Act of 2016. Public Law No. 114-255. Section 3022.  
18  
19 Available at: <https://www.congress.gov/bill/114th-congress/house-bill/34/>. Accessed  
20  
21 May 4, 2020.  
22  
23
- 24 10. Pulini AA, Caetano GM, Clautiaux H, Vergeron L, Pitts PJ, Katz G. Impact of real-world  
25  
26 data on market authorization, reimbursement decision & price negotiation [published  
27  
28 online ahead of print August 28, 2020]. *Ther Innov Regul Sci*. 2020. doi:  
29  
30 10.1007/s43441-020-00208-1  
31  
32
- 33 11. Hernandez-Boussard T, Monda KL, Crespo BC, Riskin D. Real world evidence in  
34  
35 cardiovascular medicine: ensuring data validity in electronic health record-based studies.  
36  
37 *J Am Med Inform Assoc*. 2019;26:1189-1194. doi: 10.1093/jamia/ocz119.  
38  
39
- 40 12. McCormick N, Lacaille D, Bhole V, Avina-Zubieta JA. Validity of heart failure  
41  
42 diagnoses in administrative databases: a systematic review and meta-analysis. *PLoS One*.  
43  
44 2014;9:e104519. doi: 10.1371/journal.pone.0104519.  
45  
46
- 47 13. Alqaisi F, Williams LK, Peterson EL, Lanfear DE. Comparing methods for identifying  
48  
49 patients with heart failure using electronic data sources. *BMC Health Serv Res*.  
50  
51 2009;9:237. doi: 10.1186/1472-6963-9-237  
52  
53  
54  
55  
56  
57

- 1  
2  
3 14. Xu Y, Lee S, Martin E, D'souza AG, Doktorchik CTA, Jiang J, Lee S, Eastwood CA,  
4  
5 Fine N, Hemmelgarn B, et al. Enhancing ICD-code-based case definition for heart failure  
6  
7 using electronic medical record data. *J Card Fail*. 2020;26:610-617. doi:  
8  
9 10.1016/j.cardfail.2020.04.003  
10  
11  
12 15. <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world->  
13  
14 [data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/real-world-data-assessing-electronic-health-records-and-medical-claims-data-support-regulatory)  
15  
16  
17 16. Cunningham H, Tablan V, Roberts A, Bontcheva K. Getting more out of biomedical  
18  
19 documents with GATE's full lifecycle open source text analytics. *PLoS Comput Biol*.  
20  
21 2013;9(2):e1002854.  
22  
23  
24 17. Van Rijsbergen CJ. Information Retrieval. 2nd ed. Butterworth-Heinemann. 1979.  
25  
26 18. Bozkurt B, Coats AJ, Tsutsui H, et al. Universal Definition and Classification of Heart  
27  
28 Failure: A Report of the Heart Failure Society of America, Heart Failure Association of  
29  
30 the European Society of Cardiology, Japanese Heart Failure Society and Writing  
31  
32 Committee of the Universal Definition of Heart Failure. *J Card Fail*. 2021 Mar 1:S1071-  
33  
34 9164(21)00050-6. doi: 10.1016/j.cardfail.2021.01.022.  
35  
36  
37 19. Lim YMF, Molnar M, Vaartjes I, et al. Generalizability of randomized controlled trials in  
38  
39 heart failure with reduced ejection fraction. *Eur Heart J Qual Care Clin Outcomes*.  
40  
41 2022;8(7):761-769.  
42  
43  
44 20. Quach S, Blais C, Quan H. Administrative data have high variation in validity for  
45  
46 recording heart failure. *Can J Cardiol* 2010;26:306–12.  
47  
48  
49 21. Bielinski SJ, Pathak J, Carrell DS, et al. A Robust e-Epidemiology Tool in Phenotyping  
50  
51 Heart Failure with Differentiation for Preserved and Reduced Ejection Fraction: the  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Electronic Medical Records and Genomics (eMERGE) Network. *J Cardiovasc Transl*  
4  
5 *Res.* 2015 Nov;8(8):475-83.  
6  
7  
8 22. Blecker S, Katz SD, Horwitz LI, Kuperman G, Park H, Gold A, Sontag D. Comparison of  
9  
10 approaches for heart failure case identification from electronic health record data. *JAMA*  
11  
12 *Cardiol.* 2016;1:1014-1020. doi: 10.1001/jamacardio.2016.3236  
13  
14  
15 23. Ng K, Steinhubl SR, deFilippi C, Dey S, Stewart WF. Early detection of heart failure  
16  
17 using electronic health records: practical implications for time before diagnosis, data  
18  
19 diversity, data quantity, and data density. *Circ Cardiovasc Qual Outcomes.* 2016;9:649-  
20  
21 658. doi: 10.1161/CIRCOUTCOMES.116.002797  
22  
23  
24 24. Tison GH, Chamberlain AM, Pletcher MJ, Dunlay SM, Weston SA, Killian JM, Olgin  
25  
26 JE, Roger VL. Identifying heart failure using EMR-based algorithms. *Int J Med Inform.*  
27  
28 2018;120:1-7. doi: 10.1016/j.ijmedinf.2018.09.016  
29  
30  
31 25. Khand AU, Shaw M, Gemmel I, Cleland JG. Do discharge codes underestimate  
32  
33 hospitalisation due to heart failure? Validation study of hospital discharge coding for  
34  
35 heart failure. *Eur J Heart Fail* 2005;7: 792–797.  
36  
37  
38 26. Merry AH, Boer JM, Schouten LJ, Feskens EJ, Verschuren WM, et al. Validity of  
39  
40 coronary heart diseases and heart failure based on hospital discharge and mortality data in  
41  
42 the Netherlands using the cardiovascular registry Maastricht cohort study. *Eur J*  
43  
44 *Epidemiol* 2009;24: 237–247.  
45  
46  
47 27. U.S. Food and Drug Administration. Framework for FDA's Real-World Evidence  
48  
49 Program. 2018. Available at: <https://www.fda.gov/media/120060/download>. Accessed  
50  
51 July 26, 2020.  
52  
53  
54  
55  
56  
57

- 1  
2  
3 28. Parsons A, McCullough C, Wang J, Shih S. Validity of electronic health record-derived  
4 quality measurement for performance monitoring. *J Am Med Inform Assoc.* 2012;19:604-  
5 609. doi: 10.1136/amiajnl-2011-000557  
6  
7  
8  
9  
10 29. Kao DP, Lewsey JD, Anand IS, et al. Characterization of subgroups of heart failure  
11 patients with preserved ejection fraction with possible implications for prognosis and  
12 treatment response. *Eur J Heart Fail.* 2015;17(9):925-35.  
13  
14  
15  
16 30. Uijl A, Savarese G, Vaartjes I, et al. Identification of distinct phenotypic clusters in heart  
17 failure with preserved ejection fraction. *Eur J Heart Fail.* 2021;23(6):973-982.  
18  
19  
20  
21 31. Luna D, Franco M, Plaza C, Otero C, Wassermann S, Gambarte ML, Giunta D, de Quirós  
22 FGB. Accuracy of an electronic problem list from primary care providers and specialists.  
23 *Stud Health Technol Inform.* 2013;192:417-421.  
24  
25  
26  
27 32. Singer A, Yakubovich S, Kroeker AL, Dufault B, Duarte R, Katz A. Data quality of  
28 electronic medical records in Manitoba: do problem lists accurately reflect chronic  
29 disease billing diagnoses? *J Am Med Inform Assoc.* 2016;23:1107-1112. doi:  
30 10.1093/jamia/ocw013  
31  
32  
33  
34  
35  
36  
37 33. Makam AN, Lanham HJ, Batchelor K, Samal L, Moran B, Howell-Stampley T, Kirk L,  
38 Cherukuri M, Santini N, Leykum LK, et al. Use and satisfaction with key functions of a  
39 common commercial electronic health record: a survey of primary care providers. *BMC*  
40 *Med Inform Decis Mak.* 2013;13:86. doi: 10.1186/1472-6947-13-86.  
41  
42  
43  
44  
45  
46  
47 34. Ingelsson E, Arnlov J, Sundstrom J, Lind L. The validity of a diagnosis of heart failure in  
48 a hospital discharge register. *Eur J Heart Fail* 2005;7:787-791.  
49  
50  
51  
52  
53  
54  
55  
56  
57



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

Table 1. Prespecified heart failure-specific concepts extracted from the electronic health record.

High Priority Conditions	Comorbidities	Symptoms	Findings	Medications
Congestive HF	Myocardial infarction	Angina	LVEF	Carvedilol
HF with reduced EF	Atrial fibrillation	Chest pain		Lisinopril
HF with mid-range EF	Aortic regurgitation	Dyspnea		Metoprolol
HF with preserved EF	Mitral regurgitation	Fatigue		Furosemide
	Tricuspid regurgitation	Palpitations		

HF, heart failure; EF, ejection fraction; LVEF, left ventricular ejection fraction.

Peer review only

### EHR data source and processing

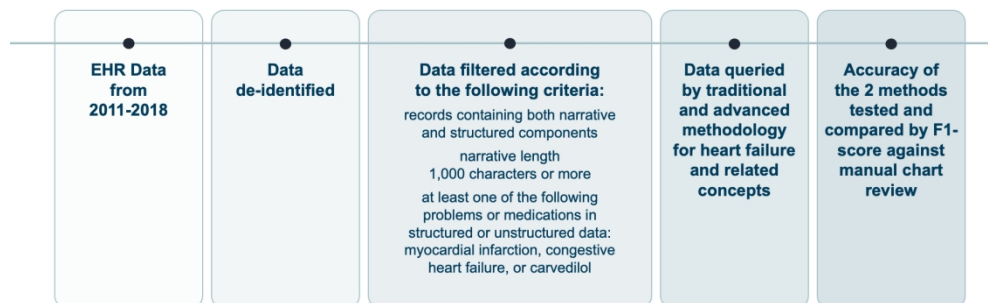


Figure 1: Electronic Health Record data source and processing.

338x190mm (144 x 144 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

### EHR data source and processing

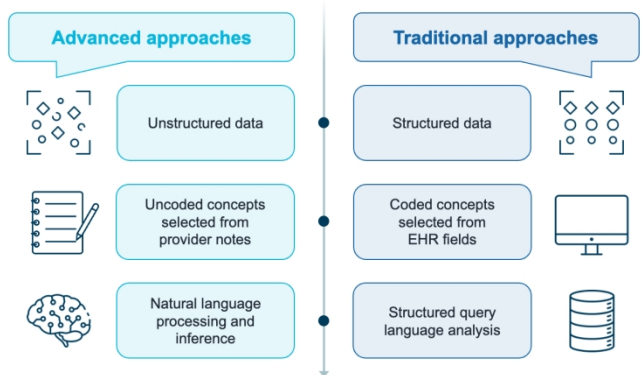


Figure 2: Comparison of traditional and advanced real-world evidence approaches. EHR, electronic health record.

338x190mm (144 x 144 DPI)

### F<sub>1</sub> score by subcategories of heart failure

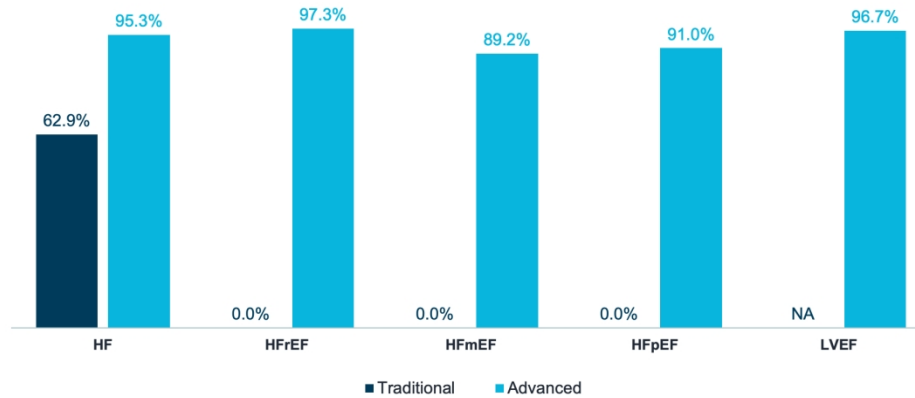
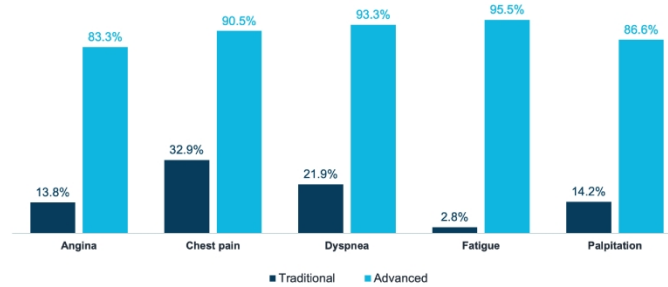
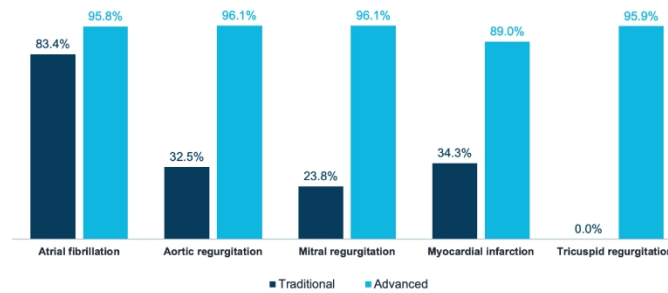
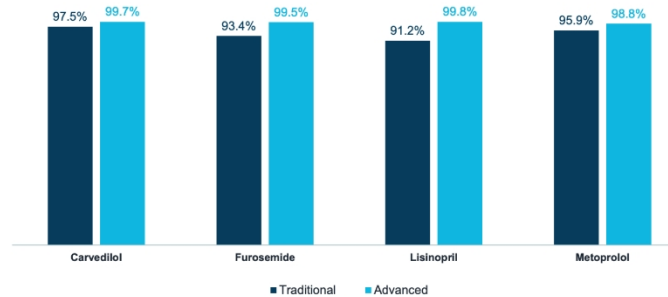


Figure 3: F1 scores for heart failure diagnoses. \*F1-score could not be calculated due to lack of data for precision. †Structured data recall is not applicable for ejection fraction because no code was available within the problem list. HF, heart failure; HFmrEF, heart failure with mildly-reduced ejection fraction; HFpEF, heart failure with preserved ejection fraction; HFrEF, heart failure with reduced ejection fraction; LVEF, left ventricular ejection fraction; 0% reflects a measured value and indicates the availability of the diagnosis code in the EHR dropdown versus N/A, not applicable, which refers to a diagnosis without available code in the relevant codeset.

338x190mm (144 x 144 DPI)

**F<sub>1</sub> score by symptoms****F<sub>1</sub> score by type of comorbidity****F<sub>1</sub> score by medication type**

F1 scores for (A) symptoms, (B) medications, and (C) comorbid conditions. \*F1 score could not be calculated due to a lack of data for precision. N/A, not applicable.

548x904mm (118 x 118 DPI)

Supplementary Table 1. Cohort identification of heart failure diagnoses, left ventricular ejection fraction, heart failure medications, symptoms, and comorbid cardiovascular conditions

	Traditional approach			Advanced approach			Concept occurrence	Encounter occurrence	P-value
	Recall, %	Precision, %	F <sub>1</sub> score, %	Recall, %	Precision, %	F <sub>1</sub> score, %			
<b>HF diagnosis</b>									
HF	46.9	95.4	62.9	96.0	94.7	95.3	265	155	<0.001
HFrEF	0	N/A*	N/A <sup>†</sup>	94.8	100.0	97.3	382	124	N/A <sup>§</sup>
HFmrEF	0	N/A*	N/A <sup>†</sup>	80.4	100.0	89.2	62	35	N/A <sup>§</sup>
HFpEF	0	N/A*	N/A <sup>†</sup>	83.5	100.0	91.0	103	55	N/A <sup>§</sup>
<b>LVEF</b>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	N/A <sup>‡</sup>	93.7	100.0	96.7	677	238	N/A <sup>§</sup>
<b>HF medications</b>									
Carvedilol	95.1	100.0	97.5	99.7	99.7	99.7	407	141	<0.001
Furosemide	87.7	100.0	93.4	99.3	99.8	99.5	1572	371	0.116
Lisinopril	83.9	100.0	91.2	99.7	99.9	99.8	1068	386	<0.001
Metoprolol	92.2	100.0	95.9	97.7	100.0	98.8	1370	397	<0.001
<b>Symptoms</b>									

Angina	7.8	60.0	13.8	84.4	82.3	83.3	265	155	<0.00 1
Chest pain	21.4	70.8	32.9	95.4	86.1	90.5	2332	756	<0.00 1
Dyspnea	12.7	78.2	21.9	94.7	92.0	93.3	4474	832	<0.00 1
Fatigue	1.4	75.0	2.8	96.5	94.5	95.5	1711	371	<0.00 1
Palpitation	8.2	52.9	14.2	90.9	82.6	86.6	896	493	<0.00 1
<b>Comorbid cardiovascular conditions</b>									
Atrial fibrillation	72.2	98.7	83.4	93.0	98.7	95.8	1214	222	<0.00 1
Aortic regurgitation	19.4	100.0	32.5	92.5	100.0	96.1	153	90	<0.00 1
Mitral regurgitation	13.5	97.1	23.8	92.8	99.6	96.1	483	185	<0.00 1
Myocardial infarction	21.1	90.9	34.3	95.5	83.4	89.0	1220	578	<0.00 1
Tricuspid regurgitation	0	N/A*	N/A†	92.2	100.0	95.9	162	78	N/A§

\*These elements did not occur when using the traditional approach. †F<sub>1</sub> scores could not be calculated due to a lack of data for precision. ‡Structured data recall is not applicable for ejection fraction because there was no code available within the problem list. §P-value could not be calculated due to the unavailability of F<sub>1</sub> scores for the traditional approach. P-values are derived from the chi-square test.

1  
2  
3 HF, heart failure; HFmrEF, heart failure with mildly reduced ejection fraction; HFpEF, heart failure with  
4 preserved ejection fraction; HFrEF, heart failure with reduced ejection fraction; LVEF, left ventricular  
5 ejection fraction; N/A, not applicable.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



## SUPPLEMENTAL MATERIAL

### Reference Standard

Traditional and advanced approaches were tested against a reference standard for physician encounters. The reference standard consisted of an independent review, with manual annotation of relevant HF-specific features, including 19 unique HF-specific concepts. For each encounter, two independent clinical annotators labeled each concept and all metadata for that concept. For example, an annotator might mark the text "DOE over last month" as dyspnea on exertion, experienced = true, current = true, relative date = 1 month. Concept occurrence was defined as the sum of all concept occurrences, allowing for multiple occurrences per encounter. Encounter occurrence was defined as the number of encounters with at least one occurrence of the concept.

Given that many concepts, such as LVEF are specific to a point in time, concepts were tested at the encounter level. For example, if a patient had an LVEF of 30% in an encounter, the data extraction would only be annotated as correct if it identified "LVEF 30%" in that specific encounter. This reference standard was used to determine accuracy of automated extracted data and structured data. Specifically, this reference standard was used to calculate recall and precision for these individual features for traditional and advanced approaches.