



Published in final edited form as:

*J Biomed Inform.* 2018 December ; 88: 62–69. doi:10.1016/j.jbi.2018.11.004.

## A Method for Harmonization of Clinical Abbreviation and Acronym Sense Inventories

Lisa V Grossman<sup>#1,2</sup>, Elliot G Mitchell<sup>#1</sup>, George Hripcsak<sup>1</sup>, Chunhua Weng<sup>1</sup>, and David K Vawdrey<sup>1,3</sup>

<sup>1</sup>Department of Biomedical Informatics, Columbia University, New York, NY, USA

<sup>2</sup>College of Physicians and Surgeons, Columbia University, New York, NY, USA

<sup>3</sup>Value Institute, NewYork-Presbyterian Hospital, New York, NY, USA

<sup>#</sup> These authors contributed equally to this work.

### Abstract

**Background:** Previous research has developed methods to construct acronym sense inventories from a single institutional corpus. Although beneficial, a sense inventory constructed from a single institutional corpus is not generalizable, because acronyms from different geographic regions and medical specialties vary greatly.

**Objective:** Develop an automated method to harmonize sense inventories from different regions and specialties towards the development of a comprehensive inventory.

**Methods:** The method involves integrating multiple source sense inventories into one centralized inventory and cross-mapping redundant entries to establish synonymy. To evaluate our method, we integrated 8 well-known source inventories into one comprehensive inventory (or *metathesaurus*). For both the metathesaurus and its sources, we evaluated the coverage of acronyms and their senses on a corpus of 1 million clinical notes. The corpus came from a different institution, region, and specialty than the source inventories.

**Results:** In the evaluation using clinical notes, the metathesaurus demonstrated an acronym (short form) microcoverage of 94.3%, representing a substantial increase over the two next largest source inventories, the UMLS LRABR (74.8%) and ADAM (68.0%). The metathesaurus demonstrated a sense (long form) micro-coverage of 99.6%, again a substantial increase compared to the UMLS LRABR (82.5%) and ADAM (55.4%).

**Conclusions:** Given the high coverage, harmonizing acronym sense inventories is a promising methodology to improve their comprehensiveness. Our method is automated, leverages the

Please address all correspondence to: Lisa V. Grossman, Columbia University Department of Biomedical Informatics, 622 W 168th St, PH-20, New York, NY 10032, United States, P: 719-244-0401, lvg2104@cumc.columbia.edu.

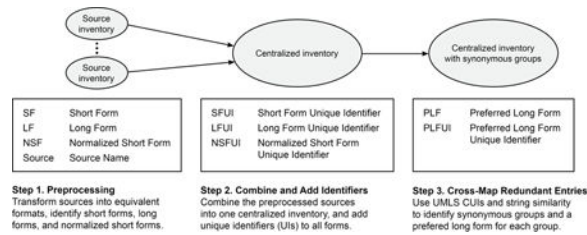
Conflicts of Interest Statement

The authors declare that they have no conflicts of interest in the research.

**Publisher's Disclaimer:** This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

extensive resources already devoted to developing institution-specific inventories in the United States, and may help generalize sense inventories to institutions who lack the resources to develop them. Future work should address quality issues in source inventories and explore additional approaches to establishing synonymy.

## Graphical abstract



## Keywords

Consumer Health Informatics; Knowledge Representation; Vocabulary and Terminology; Acronyms and Abbreviations

## 1. Introduction

In the era of consumer-driven healthcare,<sup>1</sup> more patients can instantly access their health records online than ever before.<sup>2–5</sup> Millions of patients accessed their personal health information online last year.<sup>2–5</sup> Comprehension of this information is challenging for most patients,<sup>6–10</sup> but critical to prevent medical errors,<sup>11–16</sup> increase shared decision-making,<sup>17–22</sup> and improve health outcomes.<sup>23,24</sup> According to federal policy, hospitals must allow patients to *view, download, and transmit* their own health information.<sup>25,26</sup> As a result, the percentage of healthcare organizations offering online patient portals increased from 43% in 2013 to 92% in 2015.<sup>5,27</sup> Transparency has been hailed as the next “blockbuster drug” and “healthcare revolution.” Arguably, the most promising form of transparency is giving patients access to their doctors’ notes. Over 22 million individuals accessed their doctors’ notes in early 2018, a year-over-year increase of more than 120%.<sup>28</sup> Access to notes helps patients take medications as prescribed, be better prepared for future visits, and better understand their illness conditions.<sup>29–35</sup>

However, clinical acronyms and abbreviations currently present a major barrier to patients’ understanding of their health records, especially their doctors’ notes.<sup>36,37</sup> In a previous study,<sup>37</sup> we determined that acronyms cause more misunderstanding than any other barrier, including medical terms and health literacy. At our urban academic medical center, acronyms constituted 30–50% of the words in a typical medicine admission note. In the most extreme case, a note began with an entire sentence of acronyms: ‘50 y/o f w/hx b/l SO pw/ LLQP’ (50-year-old female with a history of bilateral salpingo-oophorectomy presents with left lower quadrant pain). Therefore, it is unsurprising that acronyms cause patients to misunderstand notes. Patients’ misunderstanding may reduce the potential benefits of transparency, increase miscommunication, decrease satisfaction, increase doctors’ legal liability, and ultimately harm the patient-doctor relationship.<sup>38,39</sup> Because current interventions show extremely limited ability to improve patient’s comprehension of

acronyms,<sup>40</sup> any advances in the development of universal electronic systems to expand acronyms should have major clinical significance and far-reaching consequences such as better shared decision-making and improved health outcomes.

Electronic systems to expand acronyms rely on *sense inventories*, defined as controlled vocabularies of acronyms and their meanings (senses). Existing electronic systems currently have limited power to expand clinical acronyms,<sup>40</sup> primarily due to the lack of *comprehensiveness* (or *generalizability*) of existing sense inventories.<sup>41–44</sup> For example, sense inventories such as the Unified Medical Language System (UMLS) *LRABR* cover only 35–67% of acronyms found in doctors' notes.<sup>41,42</sup> Existing sense inventories lack generalizability because acronyms from different geographic regions and clinical specialties vary greatly.<sup>45</sup> For example, “2/2” (secondary to) is used in almost every clinical note at our medical center, but is rarely seen outside New York. “FOB” (father of baby) is frequently used in obstetrics, but not other specialties. Because of this variability, previously-developed methods have focused primarily on developing institution-specific sense inventories.<sup>45–51</sup> Institution-specific inventories, although beneficial, are labor-intensive to create and may not generalize well to clinical text from different geographical regions and specialties. Because of the labor required, developing institution-specific sense inventories at every US healthcare organization is not feasible, especially without fully automated methods which currently do not exist.

In this paper, we overcome the limitations of existing sense inventories by developing an automated methodology to harmonize sense inventories from different regions and specialties towards the development of a more comprehensive inventory. This work focuses on two critical questions: (1) can a fully automated method be developed to harmonize sense inventories? (2) does the method improve the coverage of acronyms and their senses in unseen clinical texts? Our method involves integrating multiple source sense inventories into one centralized inventory and cross-mapping redundant entries to establish synonymy. We hypothesized that a harmonized sense inventory will demonstrate greater *generalizability* than its constituent inventories. To test this hypothesis, we constructed one centralized inventory from 8 well-known source inventories, then evaluated its coverage of acronyms and their meanings (senses) in clinical notes from an institution unrelated to the source inventories. Our method leverages the extensive pre-established resources to developing institution-specific sense inventories in the United States, and may help generalize sense inventories to institutions without the resources to develop their own. A comprehensive inventory is critical for the future development of widely-available software to interpret medical acronyms for patients across the United States.

## 2. Methods

### 2.1. Method for Harmonizing Sense Inventories and Design Considerations

Harmonizing sense inventories is challenging because existing sense inventories have inconsistent formats, and may include significant redundancies. Inventories like Another Database of Abbreviations in Medline (ADAM) contain multiple synonymous expansions for single acronyms.<sup>52</sup> Sense inventories often overlap with each other in content, which creates further redundancy. Manually resolving this redundancy would be tedious and time

consuming. Instead, our proposed method seeks to *automatically resolve redundancy* by cross-mapping synonymous senses, circumventing the need for manual curation.

Figure 1 outlines the three main steps in the automated methodology for integrating sense inventories and crossmapping redundant entries. The full code and description of the method is located on GitHub.<sup>1</sup> First, source sense inventories are preprocessed to fit the same format. Second, source inventories are harmonized and non-semantic identifiers are added.<sup>53</sup> Third, synonymy between equivalent senses is automatically established using MetaMap via concept identification to form *synonymous groups*, or groups of synonymous senses. Establishing synonymy is critical to prevent redundancy and ensure concept-orientation.<sup>53–55</sup> Once generated, the final integrated sense inventory is outputted to a structured relational database designed to transparently represent source information and relationships between entries (described in Supplementary Tables 1, 2, and 3).

In designing this method, we sought to emulate successful design principles from existing endeavors to harmonize biomedical knowledge, such as the UMLS, which combines multiple biomedical terminologies into a single metathesaurus. First, the method represents all data in a common format to facilitate computational use. Second, the method preserves *source transparency*, such that source information viewed in the new format retains its original perspective and intent, and all entities can be attributed back to their source.<sup>56</sup> Third, the method uses context-free, non-semantic identifiers, which allows future evolution of the source and harmonized inventories while preserving meaning once in use.<sup>53,57</sup>

**2.1.1. Step 1. Preprocessing:** In step 1, each source database is processed to fit an equivalent format. To maintain source transparency, our method preserves all information present in each source sense inventory. During preprocessing, the short form, long form, and normalized short form for each acronym is identified from the source. The term *short form* (SF) describes the actual acronym (e.g. ‘MS’), while the term *long form* (LF) describes its spelled-out counterpart (e.g. ‘multiple sclerosis’). To create the normalized short form (NSF), the short form is converted to uppercase, punctuation is removed, and white space is removed.

**2.1.2. Step 2. Combine and Add Identifiers:** In step 2, the preprocessed source sense inventories are combined into one centralized inventory. Unique identifiers (UIs) are assigned to each unique normalized short form (NSFUI), short form (SFUI), and long form (LFUI). Each unique identifier is preceded with one letter to indicate its type (N for normalized short form, S for short form, and L for long form).

**2.1.3. Step 3. Cross-Map Redundant Entries:** In step 3, short form / long form (SF/LF) pairs with the same senses are cross-mapped into *synonymous groups* as follows. Prior to grouping, long forms are lexically normalized using UMLS lexical variant generation tools<sup>58</sup> and mapped to UMLS concepts using MetaMap.<sup>49</sup> Then, a two-step process is used to group synonymous SF/LF pairs. In the first step, long forms are grouped based on their semantic meaning. If two long forms map to the same UMLS Concept

<sup>1</sup><https://github.com/elliottgmitchell/clinical-acronym-metathesaurus>

Unique Identifier (CUI), they are identified as synonymous and grouped together. In the second step, long forms are grouped based on their character similarity, determined using the normalized Levenshtein distance metric. Levenshtein distance measures the similarity of two character strings by counting the number of insertions, deletions, and substitutions required to change one string to the other. Normalized Levenshtein distance divides the Levenshtein distance by string length to adjust for differing lengths of comparison strings.<sup>59</sup> We determined the normalized Levenshtein distance threshold for cross-mapping two entries using an iterative heuristic evaluation with 20 common abbreviations. Based on manual inspection of the distribution of normalized Levenshtein distances between normalized long forms, we selected 10 potential thresholds. Then, we manually compared the groupings for the 20 common abbreviations at each threshold, and evaluated whether synonymous long forms had grouped together without including non-synonymous terms. Based on this analysis, we identified 0.30 as the appropriate threshold for normalized Levenshtein distance.

After identifying synonymous groups, a preferred long form (PLF) is assigned to each group. The long form that is most similar to all other long forms in the synonymous group is chosen as the preferred long form, as quantified by having the smallest sum of Levenshtein distances from all other long forms in the group. Each preferred long form receives one preferred long form unique identifier (PLFUI), preceded with the letter “P.”

**2.1.4. Example:** Figure 2 provides an example to illustrate the 3-step automated method, using the acronym “DNR.” In step 1, short forms (DNR, Dnr, dnr) and long forms (daunorubicin, do not resuscitate, etc) are identified and equivalent formatting is applied. A normalized short form, “DNR,” is generated for all short forms, marking them as potentially related. In step 2, the source sense inventories are combined into one centralized inventory. Distinct short form and long form unique identifiers are generated. The normalized short form “DNR” receives only one normalized short form unique identifier. In step 3, synonymous entries such “DNR, Daunorubicin” and “DNR, Daunomycin” are cross-mapped and assigned to the same synonymous group. “Daunorubicin” and “Daunomycin” map to the same UMLS CUI, since Daunorubicin and Daunomycin refer to the same drug. A preferred long form is selected for each synonymous group and assigned an identifier.

## 2.2. Evaluation

To evaluate our method, we built a clinical abbreviation and acronym *metathesaurus* that integrates eight well-known sense inventories in the US (Table 1 and Table 2), and evaluated it on 1 million intensive care notes from the MIMIC-III corpus.<sup>60</sup> MIMIC-III is not related to any corpora used to construct the source inventories, and is also from a different geographic region and medical specialty than the source inventories.

**2.2.1. Sources:** Table 1 identifies the 8 source sense inventories incorporated into the acronym metathesaurus. (#1) The UMLS LRABR inventory, distributed by the National Library of Medicine, contains abbreviations and acronyms used to process terms for entry into the UMLS.<sup>61</sup> (#2) Another Database of Abbreviations in Medline (ADAM) contains abbreviations from titles and abstracts extracted from Medline in 2006.<sup>52</sup> (#3) Berman’s

inventory is a manually-derived database of acronyms and their long forms commonly used in pathology.

In addition to the UMLS LRABR, ADAM, and Berman, we included several sense inventories curated from institution-specific clinical corpora, which may include more clinically-oriented abbreviations. The institution-specific inventories include: (#4, #5) two inventories automatically derived from notes corpora at Vanderbilt University,<sup>62</sup> (#6) Stetson's manually curated inventory from notes at Columbia University,<sup>63</sup> and (#7) a local, manually-developed sense inventory from the Columbia University Department of Obstetrics and Gynecology (unpublished). Finally, we included Wikipedia's list of medical abbreviations (#8),<sup>64</sup> which is updated more frequently than other sources and may contain newer abbreviations.

**2.2.2. Dataset:** To perform the evaluation, we used a corpus of clinical notes from the MIMIC-III dataset.<sup>60</sup> The corpus includes over 2 million notes from over 40,000 patients who stayed in intensive care units at Beth Israel Deaconess Medical Center between 2001 and 2012. The corpus vocabulary consists of about 1.8 million words.

**2.2.3. Short Form Coverage:** Short form coverage measures whether short forms that appear in the corpus also appear in the acronym metathesaurus. We identified short forms in the first 1 million notes of MIMIC-III (hereby, the corpus) using the Clinical Abbreviation Recognition and Disambiguation (CARD) framework from Wu and colleagues.<sup>46,47</sup> For both the metathesaurus and its sources, we calculated:

Macro-coverage (Short Form)

$$= \frac{\text{Number of unique short forms in the corpus that match a short form in the metathesaurus}}{\text{Number of unique short forms in the corpus}}$$

Micro-coverage (Short Form) =

$$\frac{\text{Number of short form instances in the corpus that match a short form in the metathesaurus}}{\text{Number of short form instances in the corpus}}$$

**2.2.4. Sense Coverage:** Sense coverage measures whether the senses of short forms that appear in the corpus also appear in the acronym metathesaurus. A domain expert (LVG) manually annotated 60 random instances each for four randomly selected short forms with 1) more than one sense, and 2) more than 60 instances in the corpus. For both the metathesaurus and its sources, we calculated:

Macro-coverage (Sense) =

$$\frac{\text{Number of unique senses in the gold standard that match a short form in the metathesaurus}}{\text{Number of unique senses in the gold standard}}$$

$$\text{Micro-coverage (Sense)} = \frac{\text{Number of sense instances in the gold standard that match a short form in the metathesaurus}}{\text{Number of sense instances in the gold standard}}$$

**2.2.5. Descriptive Analysis:** We conducted descriptive analysis of the acronym metathesaurus in RStudio (R version 3.3.3), including number of senses per normalized short form and overlap between source sense inventories.

### 3. Results

#### 3.1. Evaluation

Table 3 and Figure 3 describe the acronym metathesaurus coverage and compare it with coverages for source sense inventories. On the MIMIC-III clinical notes corpus, the metathesaurus had a short form macro-coverage of 39% and a short form micro-coverage of 94%. The sense macro-coverage is 91% and the sense micro-coverage is 99.6%. The complete acronym metathesaurus is located on GitHub.<sup>2</sup>

#### 3.2. Descriptive Analysis

We discovered 376,270 SF/LF pairs in the source sense inventories, which we grouped into 105,631 synonymous groups. The metathesaurus contains 52,520 unique normalized short forms (NSFs), with an average 2.01 senses for each NSF. 11,932 NSFs possess more than one sense, and 414 NSFs possess more than 20 senses (Figure 4). The abbreviation with the most senses is “PA” (N035050), with 128 unique senses, including physician assistant, primary amenorrhea, pseudomonas aeruginosa, and psychoanalysis.

Table 4 and Figure 5 describe the degree of overlap between the source sense inventories. On average, synonymous groups contained SF/LF pairs from 1.32 source sense inventories. Out of 105,631 synonymous groups, 26,336 contained data from 2 or more sources. Table 4 provides a measure of overlap called the *group ratio*, defined as the ratio of multiple source synonymous groups to single source synonymous groups for a source. A higher group ratio indicates more overlap with other sources. Figure 5 describes how many sources the SF/LF pairs in each synonymous group come from.

### 4. Discussion

Our new automated method is promising, and is distinct from the extensive research on the semi-automated methods that generate institution-specific inventories. The new method of harmonizing sense inventories offers several advantages over existing methods. First, the new method builds on the extensive research and resources already devoted to developing institution-specific sense inventories in the United States. Second, the method can be fully automated, meaning new sense inventories can be easily incorporated as necessary. This is critical because clinical acronyms change over time, especially with new drugs and new

<sup>2</sup><https://github.com/elliottgmitchell/clinical-acronym-metathesaurus>



clinical trials. Third, the method can identify redundancies and synonymy in existing sense inventories, which may improve their quality. Fourth, the method and metathesaurus are publicly available, and the method can be used create add-ons to existing sense inventories at various healthcare institutions.

The acronym metathesaurus demonstrates high generalizability, evidenced by high coverage of acronyms in MIMIC-III. This is notable because MIMIC-III is not related to any corpora used to construct the source inventories, and is from a different geographic region and medical specialty than the source inventories. Furthermore, MIMIC-III contains real-world ICU clinical notes of various types that discuss a broad range of patient conditions with contributions from multiple specialists. As such, high coverage of acronyms in MIMIC-III indicates our method's usefulness for developing sense inventories relevant to real-world, complex natural language processing tasks. We hypothesize that our metathesaurus demonstrates high coverage because we incorporate sources with high clinical relevance, such as sense inventories derived from institutional-specific corpora. We also include Wikipedia, which potentially contains newer and more informal medical abbreviations than other sources. Unlike prior large sense inventories,<sup>45</sup> our metathesaurus does not exclude abbreviations in lower case letters.

Notably, our coverage estimates for the UMLS LRABR and ADAM differ slightly from previously reported values. Previous estimates place short form micro-coverage at 67% and sense micro-coverage at 35% for the UMLS LRABR, and 66% and 38% for ADAM.<sup>41,42</sup> In our analysis with MIMIC-III, we found 74.8% short form micro-coverage and 82.5% sense micro-coverage for the UMLS LRABR, and 68% and 55.5% for ADAM. While our short form estimates were on par with prior literature, we hypothesize that our sense coverage estimates are slightly higher because our gold standard is smaller and potentially less varied than others. In addition, while the acronym metathesaurus performed best, our estimates of short form macro-coverage were lower than expected, only 38.9% for the metathesaurus and 23.1% for the UMLS LRABR. One possible reason is the CARD framework identified candidate abbreviations which were not actually abbreviations. The framework's automation allowed us to identify short forms in all 1 million notes, at the disadvantage of identifying potential non-abbreviations such as 'lateral/', 'lip/chin', or 'rythms.' Although the short form macro-coverage of 38.9% is low, the short form micro-coverage of 94.3% is high, indicating that the metathesaurus does cover many true abbreviations.

Interestingly, the degree of overlap between source sense inventories is lower than expected. Only 3 synonymous groups of 105,631 contained entries from all 8 sources. Since overlap between sources is low, a great improvement in generalizability from integrating sources is expected. Furthermore, the low degree of overlap supports the hypothesis that source inventories generated using different methods, from different regions and specialties, contain different acronyms. Interestingly, the group ratio, which measures overlap between one source and all other sources, does not appear to correlate with sense coverage. This may be because a high degree of overlap could have two meanings. First, high overlap may indicate the source is not comprehensive, as all its entries occur in other sources. Second, high overlap may indicate the source is very comprehensive, as it contains entries found in many other sources.



Given the rapid adoption of electronic health records over the past ten years, an increasing need exists to interpret abbreviations and acronyms. While the use case of patient-facing applications drove the development of the acronym metathesaurus, additional applications include clinical decision support tools, tools to support interoperability across electronic health records, and teaching tools for medical students and residents. Because many abbreviations and acronyms possess more than one sense, abbreviation sense disambiguation is an important task in clinical natural language processing. Abbreviation sense disambiguation relies heavily on a complete and consistent sense inventory.<sup>43</sup>

Current semi-automated methods that generate sense inventories for a specific corpus, such as CARD, will likely perform better for acronym sense disambiguation on that specific corpus. CARD uses automated extraction and manual annotation to identify the most relevant senses for a given corpus. However, our method may perform better for acronym sense disambiguation on unknown or multiple corpora, as in the use case presented in the Introduction. While our method demonstrated impressive coverage, the volume of possible senses may add complexity to abbreviation sense disambiguation tasks. More work is necessary to determine how harmonized sense inventories perform for disambiguation tasks. We hope the open-source and comprehensive nature of the metathesaurus will facilitate such work.

#### 4.1. Limitations

The acronym metathesaurus may not perform as well on clinical corpora from a single specialty or institution compared to an inventory developed on a corpus from that specialty or institution. Although our results demonstrate that our method is valid and may solve the problem of generating comprehensive sense inventories, challenges still exist. Currently, the method does not address quality issues in source inventories. Future work should also explore methods for improving the quality of source sense inventories at the time of integration. Furthermore, MetaMap only identified concepts in about 30% of entries, even after lexical normalization and optimization, suggesting that clinical named entity recognition tools like MetaMap alone are insufficient to establish synonymy between all entries. While the Levenshtein distance cutoff enables automatic creation of synonymous groups, the cutoff may not properly group all synonymous long forms. Future work should explore possible additional methods for assigning preferred long form and creating the synonymous groups, including ways to increase MetaMap performance, and evaluate the quality of these methods.

## 5. Conclusion

We developed an automated method to harmonize sense inventories from different regions and specialties towards the development of a comprehensive inventory. In an evaluation using clinical notes, an acronym metathesaurus constructed from 8 source sense inventories demonstrated an acronym (short form) coverage of 94.3% and a sense (long form) coverage of 99.6%, representing an increase in coverage over well-known sense inventories such as the UMLS LRABR and ADAM. Harmonizing sense inventories is therefore a promising methodology to improve their comprehensiveness.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

This work was supported by the National Library of Medicine (T15LM007079, R01LM006910, PI: Hripcsak; R01LM009886, PI: Weng) and the Agency for Healthcare Research and Quality (R01HS21816, PI: Vawdrey).

## 7. References

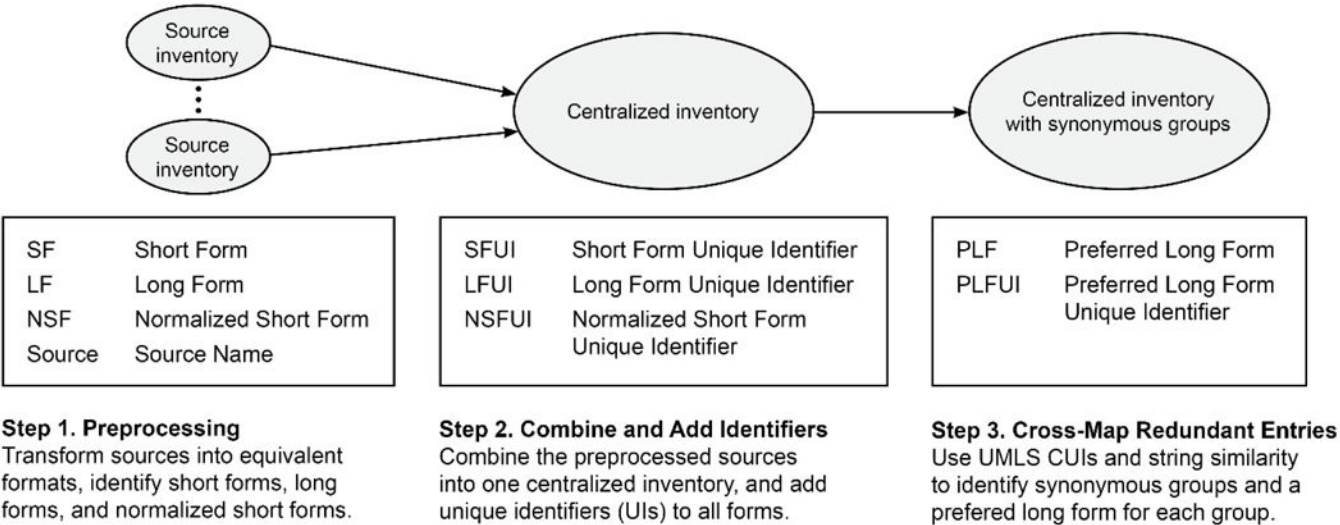
1. Berwick D Era 3 for Medicine and Health Care. JAMA vol: 315, 1329–30 (2016). [PubMed: 26940610]
2. Delbanco T et al. Open notes: doctors and patients signing on. Ann. Intern. Med. 153, 121–5 (2010). [PubMed: 20643992]
3. Delbanco T, Walker J, Bell SK, et al. Inviting Patients to Read Their Doctors' Notes: A Quasi-experimental Study and a Look Ahead. Ann Intern Med. 2012;157(7):461. [PubMed: 23027317]
4. Walker J, Leveille SG, Ngo L, et al. Inviting patients to read their doctors' notes: patients and doctors look ahead: patient and physician surveys. Ann Intern Med. 2011;155(12):811–9. [PubMed: 22184688]
5. American Hospital Association. Individuals' Ability to Electronically Access Their Hospital Medical Records, Perform Key Tasks is Growing. (2016).
6. Irizarry T et al. Patient Portals as a Tool for Health Care Engagement: A Mixed-Method Study of Older Adults With Varying Levels of Health Literacy and Prior Patient Portal Use. J. Med. Internet Res. 19, e99 (2017). [PubMed: 28360022]
7. Irizarry T, De Vito Dabbs A & Curran CR Patient portals and patient engagement: A state of the science review. J. Med. Internet Res. 17, e148 (2015). [PubMed: 26104044]
8. Health Literacy: A Prescription to End Confusion Institute of Medicine (US) Committee on Health Literacy; Nielsen-Bohlman L, Panzer AM, Kindig DA, editors. Washington (DC): National Academies Press (US); 2004
9. Sarkar U, Karter AJ & Liu JY The Literacy Divide: Health Literacy and the Use of an Internet-Based Patient Portal in an Integrated Health System—Results from the Diabetes Study of Northern California (DISTANCE). J Heal. Commun 15, 183–196 (2010).
10. Goel MS et al. Patient reported barriers to enrolling in a patient portal. J. Am. Med. Informatics Assoc. 18, i8–i12 (2011).
11. Weingart SN et al. Lessons from a patient partnership intervention to prevent adverse drug events. Int. J. Qual. Heal. Care 16, 499–507 (2004).
12. Weingart SN et al. Medication safety messages for patients via the web portal: The MedCheck intervention. Int. J. Med. Inform. 77, 161–168 (2008). [PubMed: 17581772]
13. Heyworth L et al. Engaging patients in medication reconciliation via a patient portal following hospital discharge. J. Am. Med. Informatics Assoc. 21, e157–e162 (2014).
14. Schnipper JL et al. Effects of an online personal health record on medication accuracy and safety: a cluster-randomized trial. J. Am. Med. Informatics Assoc. 19, 728–734 (2012).
15. Dullabh P, Sondheimer N, Katsh E, Evans MA. How patients can improve the accuracy of their medical records. eGEMS, 2 (3) (2014), p. 1080 [PubMed: 25848614]
16. Staroselsky M et al. Improving electronic health record (EHR) accuracy and increasing compliance with health maintenance clinical guidelines through patient access and input. Int. J. Med. Inform. 75, 693–700 (2006). [PubMed: 16338169]
17. Caligtan CA, Carroll DL, Hurley AC, Gersh-Zaremski R & Dykes PC Bedside information technology to support patient-centered care. Int. J. Med. Inform. 81, 442–451 (2012). [PubMed: 22285034]
18. Dalal AK et al. A web-based, patient-centered toolkit to engage patients and caregivers in the acute care setting: A preliminary evaluation. J. Am. Med. Informatics Assoc. 23, 80–87 (2016).

19. Stade D & Dykes P Nursing Leadership in Development and Implementation of a Patient-Centered Plan of Care Toolkit in the Acute Care Setting. *CIN Comput. Informatics, Nurs.* 33, 90–92 (2015).
20. Maher M et al. A Novel Health Information Technology Communication System to Increase Caregiver Activation in the Context of Hospital-Based Pediatric Hematopoietic Cell Transplantation: A Pilot Study. *JMIR Res. Protoc.* 4, e119 (2015). [PubMed: 26508379]
21. Maher M et al. User-Centered Design Groups to Engage Patients and Caregivers with a Personalized Health Information Technology Tool. *Biol. Blood Marrow Transplant.* 22, 349–358 (2016). [PubMed: 26343948]
22. Dykes PC et al. Building and testing a patient-centric electronic bedside communication center. *J. Gerontol. Nurs.* 39, 15–9 (2013).
23. Kruse CS, Bolton K & Freriks G The effect of patient portals on quality outcomes and its implications to meaningful use: A systematic review. *J. Med. Internet Res.* 17, 1–8 (2015).
24. Mold F et al. Patients' online access to their electronic health records and linked online services: A systematic review in primary care. *Br. J. Gen. Pract.* 65, e141–e151 (2015). [PubMed: 25733435]
25. Bitton A, Poku M & Bates D in *Information Technology for Patient Empowerment in Healthcare* (eds. Grando M, Rozenblum R & Bates D) 75–90. (Walter de Gruyter Inc, 2015).
26. 2014 Edition EHR Certification Criteria Grid Mapped to Meaningful Use Stage 2. Available at: [https://www.healthit.gov/sites/default/files/2014editionehrcertificationcriteria\\_mustage2.pdf](https://www.healthit.gov/sites/default/files/2014editionehrcertificationcriteria_mustage2.pdf).
27. HealthIT.gov National Learning Consortium. How to Optimize Patient Portals for Patient Engagement and Meet Meaningful Use Requirements. (2013).
28. [OpenNotes.org](https://www.opennotes.org/) (2018). Available at: <https://www.opennotes.org/>.
29. Wolff JL et al. Inviting patients and care partners to read doctors' notes: OpenNotes and shared access to electronic medical records. *J. Am. Med. Inform. Assoc.* 157, 461–470 (2016).
30. Bell SK et al. A patient feedback reporting tool for OpenNotes: Implications for patient-clinician safety and quality partnerships. *BMJ Qual. Saf.* 26, 312–322 (2017).
31. Nazi KM, Turvey CL, Klein DM, Hogan TP & Woods SS VA OpenNotes: exploring the experiences of early patient adopters with access to clinical notes. *J. Am. Med. Informatics Assoc.* 22, 380–389 (2014).
32. Leveille SG et al. Evaluating the impact of patients' online access to doctors' visit notes: designing and executing the OpenNotes project. *BMC Med. Inform. Decis. Mak.* 12, 32 (2012). [PubMed: 22500560]
33. Bell SK, Mejilla R, Anselmo M et al. When doctors share visit notes with patients: a study of patient and doctor perceptions of documentation errors, safety opportunities and the patient–doctor relationship. *BMJ Qual Saf.* 2017;26(4):262–270.
34. Gerard M, Fossa A, Folcarelli PH, Walker J & Bell SK What patients value about reading visit notes: A qualitative inquiry of patient experiences with their health information. *J. Med. Internet Res.* 19, (2017).
35. Goldzweig CL Pushing the Envelope of Electronic Patient Portals to Engage Patients in Their Care. *Ann. Intern. Med.* 157, 525 (2012). [PubMed: 23027322]
36. Keselman A, et al. Towards consumer-friendly PHRs: patients' experience with reviewing their health records, in: *AMIA Annu. Symp. Proc 2007*, pp. 399–403.
37. Grossman L, Masterson Creber R, Restaino S & Vawdrey DK. Sharing Clinical Notes with Hospitalized Patients via an Acute Care Portal. *AMIA Annu. Symp. Proc* (2017), pp. 800–809.
38. Manson A Language concordance as a determinant of patient compliance and emergency room use in patients with asthma. *Med. Care* 26, 1119–1128 (1988). [PubMed: 3199910]
39. Waitzkin H Doctor-Patient Communication. *JAMA* 252, 2441 (1984). [PubMed: 6481931]
40. Ramesh BP, Houston T, Brandt C, Fang H & Yu H Improving patients' electronic health record comprehension with NoteAid. in *Studies in Health Technology and Informatics* 192, 714–718 (2013). [PubMed: 23920650]
41. Liu H, Lussier YA & Friedman C A study of abbreviations in the UMLS. *Proceedings. AMIA Annu. Symp.* Proc 393–7 (2001).
42. Xu H, Stetson PD & Friedman C A study of abbreviations in clinical notes. *AMIA Annu. Symp. Proc* 821–5 (2007).

43. Moon S, McInnes B & Melton GB Challenges and practical approaches with word sense disambiguation of acronyms and abbreviations in the clinical domain. *Healthc. Inform. Res.* 21, 35–42 (2015). [PubMed: 25705556]
44. Wu Y. A comparative study of current Clinical Natural Language Processing systems on handling abbreviations in discharge summaries; *AMIA Annu. Symp. Proc.*; 2012. 997–1003.
45. Moon S, Pakhomov S, Liu N, Ryan JO & Melton GB A sense inventory for clinical abbreviations and acronyms created using clinical notes and medical dictionary resources. *J. Am. Med. Informatics Assoc.* 21, 299–307 (2014).
46. Wu, et al. A long journey to short abbreviations: developing an open-source framework for clinical abbreviation recognition and disambiguation (CARD), *Journal of the American Medical Informatics Association*, Volume 24, Issue e1, 1 4 2017, Pages e79–e86. [PubMed: 27539197]
47. Xu H, Stetson PD & Friedman C Methods for Building Sense Inventories of Abbreviations in Clinical Notes. *J. Am. Med. Informatics Assoc.* 16, 103–108 (2009).
48. Dannélls D. Automatic acronym recognition; *Proc. Elev. Conf. Eur. Chapter Assoc. Comput. Linguist. Posters Demonstr. - EACL*; 2006. 167
49. MetaMap - A Tool For Recognizing UMLS Concepts in Text. Available at: <https://metamap.nlm.nih.gov/>. (Accessed: 31st March 2016)
50. Wu Y et al. A Preliminary Study of Clinical Abbreviation Disambiguation in Real Time. *Appl. Clin. Inform.* 6, 364–74 (2015). [PubMed: 26171081]
51. Wu Y et al. Clinical acronym/abbreviation normalization using a hybrid approach. *CEUR Workshop Proc.* 1179, (2013).
52. Zhou W, Torvik VI & Smalheiser NR ADAM: Another database of abbreviations in MEDLINE. *Bioinformatics* 22, 2813–2818 (2006). [PubMed: 16982707]
53. Cimino JJ Desiderata for controlled medical vocabularies in the twenty-first century. *Methods Inf. Med.* 37, 394–403 (1998). [PubMed: 9865037]
54. Bodenreider O The Unified Medical Language System (UMLS): integrating biomedical terminology. *Nucleic Acids Res.* 32, D267–D270 (2004). [PubMed: 14681409]
55. Cimino JJ Auditing the Unified Medical Language System with Semantic Methods. *J. Am. Med. Informatics Assoc.* 5, 41–51 (1998).
56. Hole WT et al. Achieving ‘source transparency’ in the UMLS Metathesaurus. *Stud. Health Technol. Inform.* 107, 371–5 (2004). [PubMed: 15360837]
57. Cimino JJ In defense of the Desiderata. *J. Biomed. Inform.* 39, 299–306 (2006). [PubMed: 16386470]
58. Lexical Tools, 2017 Release. Available at: <https://lsg3.nlm.nih.gov/LexSysGroup/Projects/lvg/%0Acurrent/web/index.html%0A>. (Accessed: 31st March 2016)
59. Levenshtein VI Binary codes capable of correcting deletions, insertions, and reversals. *Sov. Phys. Dokl.* 10, 707–710 (1966).
60. Johnson AEW et al. MIMIC-III, a freely accessible critical care database. *Sci. data* 3, 160035 (2016). [PubMed: 27219127]
61. UMLS Reference Manual. Available at: <https://www.ncbi.nlm.nih.gov/books/NBK9680/>. (Accessed: 31st March 2016)
62. Recognition and Disambiguation of Clinical Abbreviations. Available at: <https://sbmi.uth.edu/ccb/resources/abbreviation.htm>. (Accessed: 12th March 2016)
63. Stetson PD, Scotch M & Hripcsak G The sublanguage of cross-coverage. *AMIA Annu. Symp. Proc* 742–6 (2002).
64. Wikipedia: List of Medical Abbreviations. Available at: [https://en.wikipedia.org/wiki/List\\_of\\_medical\\_abbreviations](https://en.wikipedia.org/wiki/List_of_medical_abbreviations). (Accessed: 12th March 2016)
65. Berman JJ Pathology Abbreviated: A Long Review of Short Terms. *Arch. Pathol. Lab. Med.* 128, 347–352 (2004). [PubMed: 14987146]

**Highlights**

- We devised a method to harmonize acronym and abbreviation sense inventories
- To evaluate our method, we harmonized 8 existing inventories into one metathesaurus
- The harmonized sense inventory had a higher coverage of acronyms in clinical texts
- The harmonized sense inventory was more comprehensive than existing inventories
- The harmonized sense inventory generalized to another institution's clinical texts



**Figure 1.**  
Automated Method for Harmonizing Sense Inventories



## Step 1. Preprocessing

ADAM	Short Form	Long Form
	DNR	do-not-resuscitate
	DNR	do not resuscitate
	DNR	not resuscitate
	DNR	daunorubicin
	DNR	daunomycin
Berman	Short Form	Long Form
	dnr	do not resuscitate
	dnr	daunorubicin
UMLS LRABR	Short Form	Long Form
	DNR	digital noise reduction
	DNR	do-not-resuscitate order
	DNR	do-not-resuscitate
	DNR	do not resuscitate order
	DNR	do not resuscitate
	DNR	daunorubicin
	DNR	diffusion-to-noise ratio
	DNR	daunorubicin
	DNR	do-not-resuscitate
	DNR	dose nonuniformity ratio
	DNR	dose non-uniformity ratio
	Dnr	daunorubicin
	dnr	daunorubicin
	dnr	do not resuscitate
UT Health	Short Form	Long Form
	dnr	do not resuscitate
Vanderbilt	Short Form	Long Form
	DNR	do not resuscitate
Wikipedia	Short Form	Long Form
	DNR	do not resuscitate

## Step 2. Combine and Add Identifiers

Normalized Short Form: DNR [N013686]			
SFUI	SF	LFUI	LF
S018267	DNR	L063893	do-not-resuscitate
S018267	DNR	L000003	do not resuscitate
S018267	DNR	L109767	not resuscitate
S018267	DNR	L059961	daunorubicin
S018267	DNR	L059958	daunomycin
S018267	DNR	L062493	digital noise reduction
S018267	DNR	L063894	do-not-resuscitate order
S018267	DNR	L063891	do not resuscitate order
S018267	DNR	L063890	do not resuscitate
S018267	DNR	L062417	diffusion-to-noise ratio
S018267	DNR	L064285	dose nonuniformity ratio
S018267	DNR	L064284	dose non-uniformity ratio
S019150	Dnr	L059961	daunorubicin
S079540	dnr	L059961	daunorubicin
S079540	dnr	L063890	do not resuscitate

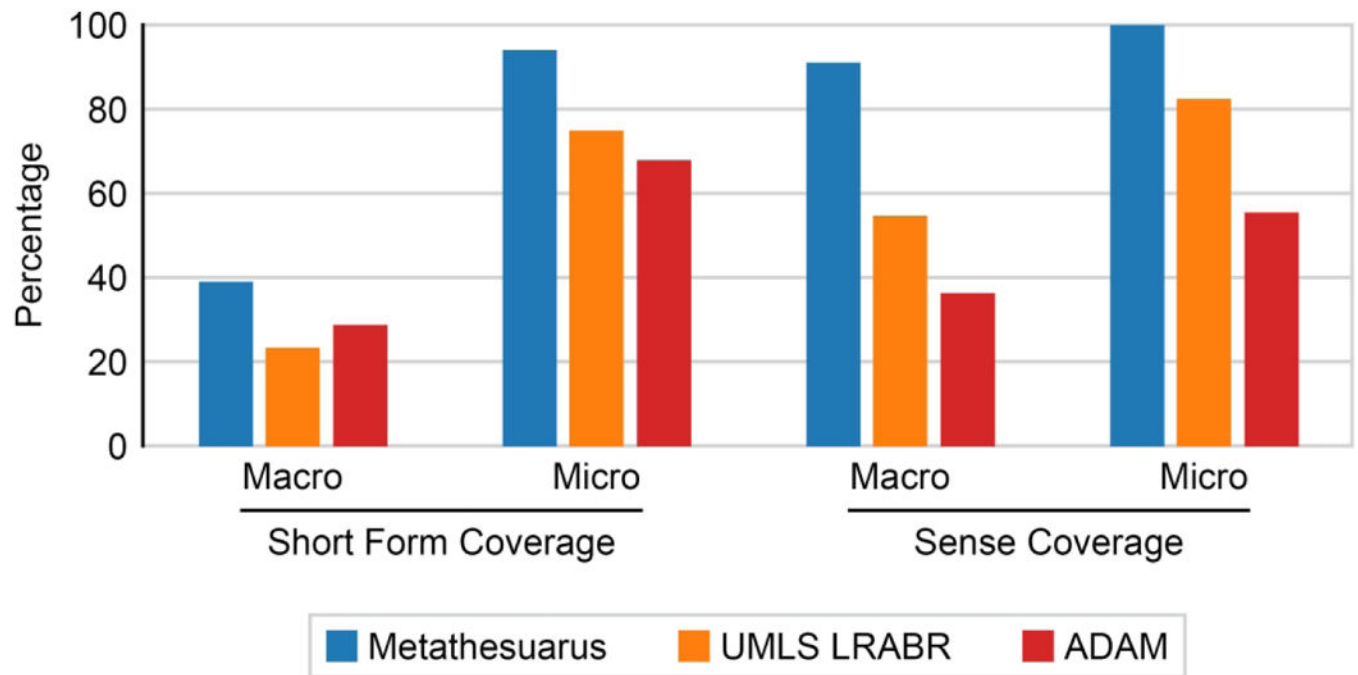
## Step 3. Cross-Map Redundant Entries

Normalized Short Form: DNR [N013686]	
PLFUI	PLF
P000004	daunorubicin
P000005	do not resuscitate
P000006	diffusion-to-noise ratio
P000007	digital noise reduction
P000008	dose nonuniformity ratio

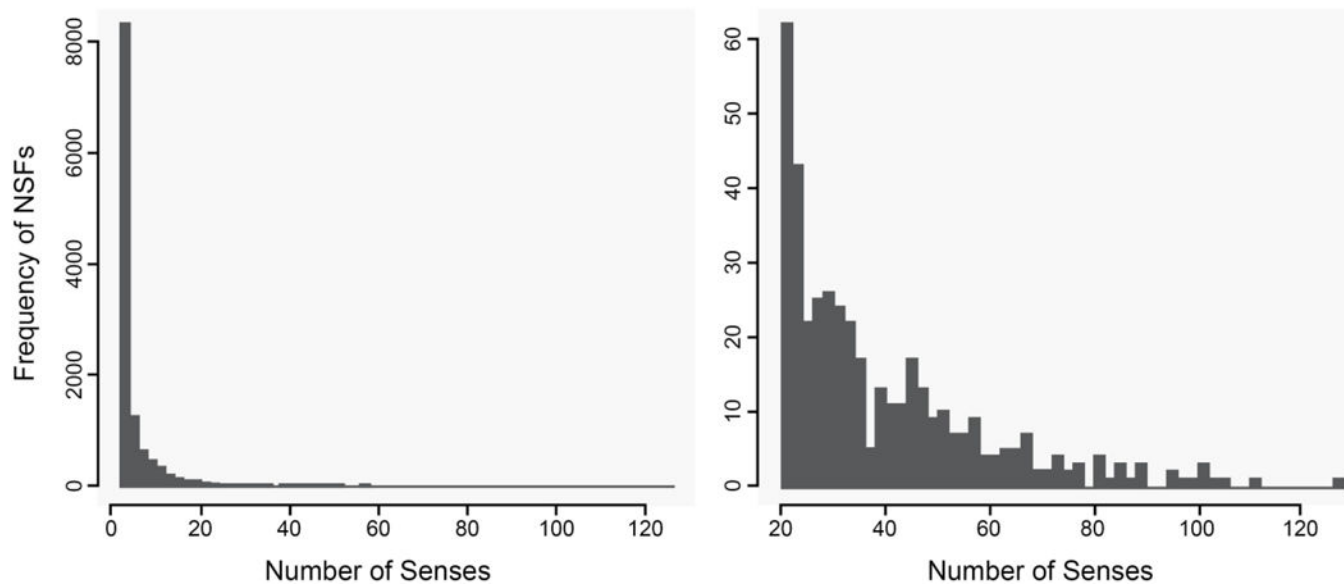
**Figure 2.**

Example Harmonization of the Acronym “DNR”

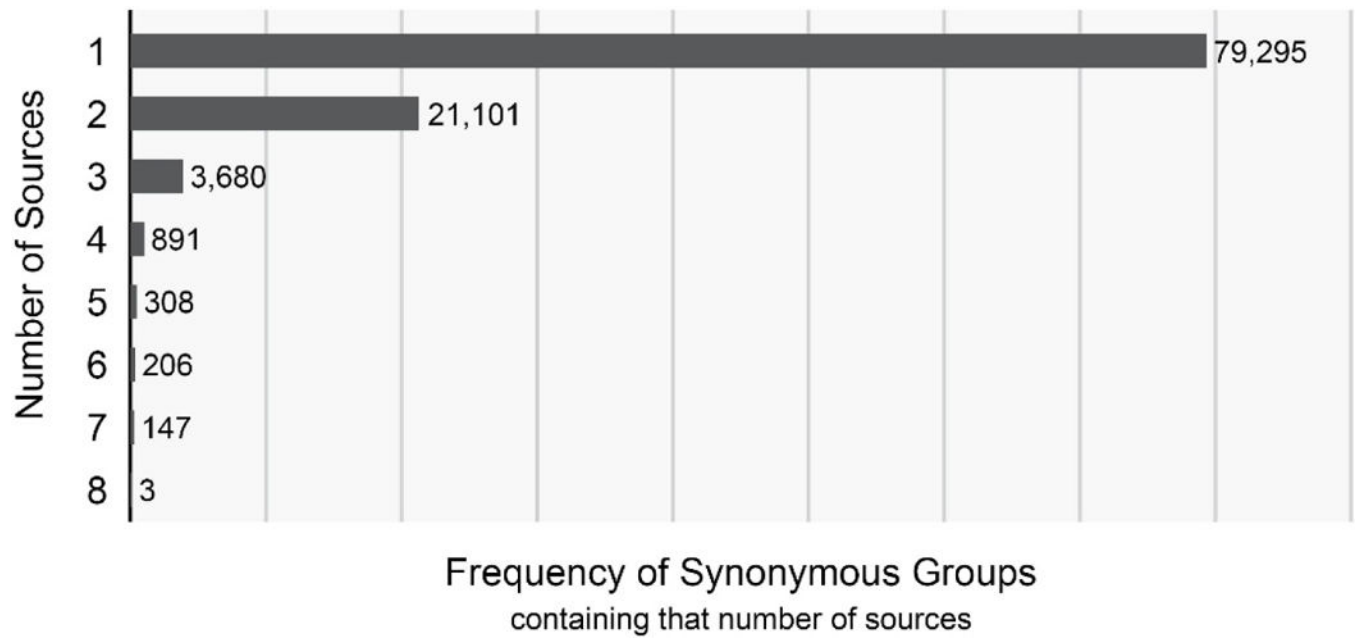
Abbreviations: SF, Short Form; LF, Long Form; SFUI, Short Form Unique Identifier; LFUI, Long Form Unique Identifier; PLF, Preferred Long Form; PLFUI, Preferred Long Form Unique Identifier.



**Figure 3.**  
Coverage of the Metathesaurus, UMLS LRABR, and ADAM



**Figure 4.**  
Number of Senses Per Normalized Short Form



**Figure 5.**  
Number of Sources Represented Within Synonymous Groups

**Table 1.**

## Source Sense Inventories Incorporated into the Metathesaurus

No.	Source
1	Unified Medical Language System (UMLS) abbreviations and acronyms inventory, LRABR <sup>61</sup>
2	Another Database of Abbreviations in Medline (ADAM) <sup>52</sup>
3	Berman's 12000 pathology abbreviations <sup>65</sup>
4	Sense inventories derived from a corpus of discharge notes at Vanderbilt University <sup>62</sup>
5	Sense inventories derived from a corpus of sign-out notes at Vanderbilt University <sup>62</sup>
6	Stetson's manually-curated sense inventory from sign-out notes at Columbia University <sup>63</sup>
7	A locally-developed sense inventory from obstetrics and gynecology notes at Columbia University
8	Wikipedia's list of medical abbreviations <sup>64</sup>

**Table 2.**

Statistics for Source Sense Inventories

No.	Source	Last updated	SF/LF pairs	Unique SFs	Unique LFs
1	UMLS LRABR	2018	261389	64370	117615
2	ADAM	2007	94657	42465	54057
3	Berman	2004	12084	7546	11156
4	Vanderbilt Discharge Summaries	2013	2090	1690	1281
5	Vanderbilt Clinical Notes	2013	2414	1929	1412
6	Stetson	2002	765	448	671
7	Columbia OBGYN	2018	219	217	212
8	Wikipedia	2018	2652	2259	2523

Abbreviations: SF, short form; LF, long form.



**Table 3.**

Coverage of the Metathesaurus and its Source Sense Inventories

Sense Inventory	<u>Short Form Coverage (%)</u>		<u>Sense Coverage (%)</u>	
	Macro	Micro	Macro	Micro
Metathesaurus	38.9	94.3	90.9	99.6
UMLS LRABR	23.1	74.8	54.5	82.5
ADAM	28.6	68.0	36.4	55.4
Berman	8.9	25.4	45.5	66.3
Vanderbilt Discharge Summaries	5.8	67.3	72.7	91.3
Vanderbilt Clinical Notes	6.4	67.9	72.7	90.8
Stetson	1.6	19.2	63.6	74.6
Columbia OBGYN	0.7	1.8	0.0	0.0
Wikipedia	6.7	50.2	27.3	53.3

**Table 4.**

## Group Ratio for Source Sense Inventories

Source Sense Inventory	Group Ratio
UMLS LRABR	0.521
ADAM	0.900
Berman	2.012
Stetson	2.621
Columbia OBGYN	2.944
Wikipedia	3.301
Vanderbilt Clinical Notes	3.497
Vanderbilt Discharge Summaries	3.611