

AcX: system, techniques, and experiments for Acronym eXpansion [Extended Version]

João L. M. Pereira
INESC-ID, IST Universidade de Lisboa, and
University of Amsterdam
joaoplmpereira@tecnico.ulisboa.pt

Helena Galhardas
INESC-ID and IST, Universidade de Lisboa
helenagalhardas@tecnico.ulisboa.pt

João Casanova*
Hitachi Vantara
joao.casanova@hitachivantara.com

Dennis Shasha
Courant Institute, New York University
shasha@cs.nyu.edu

Abstract

In this information-accumulating world, each of us must learn continuously. To participate in a new field, or even a sub-field, one must be aware of the terminology including the acronyms that specialists know so well, but newcomers do not.

Building on state-of-the-art acronym tools, our *end-to-end acronym expander system* called AcX takes a document, identifies its acronyms, and suggests expansions that are either found in the document or appropriate given the subject matter of the document. As far as we know, AcX is the first open source and extensible system for acronym expansion that allows mixing and matching of different inference modules. As of now, AcX works for English, French, and Portuguese with other languages in progress.

This paper describes the design and implementation of AcX, proposes *three new acronym expansion benchmarks*, compares state-of-the-art techniques on them, and proposes ensemble techniques that improve on any single technique. Finally, the paper evaluates the performance of AcX in end-to-end experiments on a human-annotated dataset of Wikipedia documents. Our experiments show that human performance is still better than the best automated approaches. Thus, achieving Acronym Expansion at a human level is still a rich and open challenge.

1 Introduction

Take a great historical literary figure of any culture and put him or her in the present. That person might barely understand a newspaper headline partly because of the acronyms that have no expansion. For example, a headline of the Washington Post on May 6, 2022 reads: "FDA limits use of J&J vaccine over rare blood clots". Even Shakespeare would be at a loss.

Contemporary scholars encounter similar challenges when entering a new field or a sister field. A typical document about wireless communication is practically unintelligible to a computer scientist with its talk of 3gPP, 5G, BS etc. Documents written for specialists often neglect even to define the acronyms they use [6].

Further, the proper expansion of an acronym depends on context. For example, "ISBN" can mean International Standard Book Number in a publishing context, Integrated Satellite Business Network in a satellite communication context, and International Society for

Behavioral Neuroscience in a cognitive scientific context. Thus, any system that hopes to help readers understand the intended meaning of an undefined acronym in an document must expand that acronym using its context.

1.1 High Level Architecture of an Acronym Expansion System

An *end-to-end acronym expander system* comprises the following two components:

- (i) Extraction of each acronym and (when present) its expansion within a text. For example, if a given text has "ISBN (Integrated Satellite Business Network)" then "ISBN" would be the acronym and "Integrated Satellite Business Network" would be the expansion. We call this *in-expansion* because this can be done for a particular document based solely on its own text.
- (ii) In the case that an acronym is not expanded in the text of a document, *out-expansion* chooses an expansion from a large parsed corpus (training corpus) of other documents (e.g., Wikipedia).

This paper makes the following system and data contributions:

- The **end-to-end Acronym eXpander (AcX) system** accepts a text document as input and outputs a list of acronym-expansion pairs for the acronyms found in the document, whether or not the expansions are in the document. As far as we know, AcX is the first open source and extensible system for acronym expansion that allows both the mixing and the combination of different inference modules. AcX is easily extendible to other languages. We already have versions in English, French and Portuguese. Students with no previous background in natural language processing or machine learning have already done so for French and Portuguese.
- A **benchmark of in-expansion techniques (in-expansion benchmark)**. We make use of four biomedical datasets previously proposed in the literature (i.e., Medstract [29], BIOADI [29], Schwartz and Hearst [29], and Ab3p [29]) and one of the biggest sentence based datasets from the scientific domain (i.e., SciAI [76]). Additionally, we have created a new dataset composed of Wikipedia documents from the *Computing* category. We calculate the precision, recall, and F1-measure for the extraction of both acronyms and acronym-expansion pairs. In addition, we measure training

*This work was performed while the author was a MSc student at IST, Universidade de Lisboa.

and execution times. Additionally, we have implemented, integrated and evaluated two rule-based (Schwartz and Hearst [63], MadDog [74]) and two machine-learning (based on SciBERT [7] used in [68], and SciDr [68]) acronym-expansion extraction techniques. Our experimental results show that there is not a single best technique for every dataset. Moreover, rule-based techniques overall achieve better precision while machine-learning techniques achieve better recall.

- **A benchmark of out-expansion techniques (out-expansion benchmark).** We evaluate out-expansion techniques on three datasets from different domains previously used in related work that contain documents (i.e., MSH [58], SciWISE [58], and CS Wiki [73]) and one that is constructed from independent sentences from the scientific domain (i.e., SciAD [75] revised by Egan and Bohannon [21]).

Some of the out-expansion techniques we consider have already been proposed in related work: Classic Context Vector [2, 42, 58], Surrounding Based Embedding [42], Thakker et al. [73], and Unsupervised Abbreviation Disambiguation [16]. Most recently, competitors of the SDU@AAAI competition [75] mainly used pre-trained language models based on Transformer neural networks like BERT [20] and SciBERT [7]. For purposes of out-expansion, we have adapted and evaluated: (i) the out-expander of the MadDog system [74]; (ii) SciDr [68] (which was originally hard coded to apply to the sentences of the SDU@AAAI competition, but we have extended it to multiple datasets of documents); (iii) LUKE (Language Understanding with Knowledge-based Embeddings) [81] originally for Entity Disambiguation, but we extended it to use acronyms and expansions; and (iv) techniques from Natural Language Processing (i.e., Term Frequency–Inverse Document Frequency [35], Latent Dirichlet Allocation [9], Doc2Vec [41], and Sentence Bidirectional Encoder Representations from Transformers (SBERT) [61]). We have also embedded the outputs of a variety of out-expansion techniques as features for machine learning classifiers. The net result is that AcX is an extensible system able to create a new set of out-expansion pipelines out of combinations of user-chosen techniques.

We evaluate the techniques using out-expansion accuracy (most common for these kind of disambiguation problems) and the macro F1-measure (used in the SDU@AAAI competition). We also measure execution times for both technique training and document processing. Our results show that ensembler techniques give the best accuracies and macro F1-measure. Cossim with SBERT and SciDr out-expander are the individual techniques that score best. At a slight loss in accuracy, Classic Context Vectors, or even Cossim or SVMs with Doc2Vec are almost 10 times fast on average.

- **A benchmark of end-to-end acronym expander systems (end-to-end benchmark).** We create the first end-to-end dataset of human-annotated documents that includes both in- and out-expansions. We have built a human-generated end-to-end benchmark because previous annotated datasets used automatic mechanisms to identify acronyms and those automatic techniques are neither accurate nor complete. Thus, human annotation offers a kind of gold

standard. We use this dataset to test AcX as well as the state-of-the-art acronym expansion systems. We use all English Wikipedia documents to train it. We compare the MadDog system [74] against different pipelines of AcX. We also compare AcX’s performance to that of human annotators. Those systems are evaluated using the precision, recall, and F1-measure for the acronym-expansion pairs returned by a system for each document. We also measure execution times. Our best system pipeline correctly expands most acronyms (54.97% of F1-Measure, only about 27% worse than humans) taking on average less than 2s per document.

- **A link follower component** within AcX system that infers expansions based on the links of unexpanded acronyms in input documents. Link following improves the F1-measure of the best AcX pipeline identified by our end-to-end benchmark. Our results show that our best system pipeline slight improves, +0.37 of F1-measure, when this component is used.

This paper is organized as follows: Section 2 presents related work, particularly techniques for in-expansion extraction and techniques for acronym out-expansion, including references to entity linking. Section 3 describes the AcX system and its components. The next three sections (Sections 4, 5 6) describe the proposed benchmarks and the corresponding analyses of the results obtained from the experiments performed in the context of each benchmark. Finally, Section 7 presents the main conclusions and ideas for future work.

2 Related Work

This section describes the work that is closely relevant to acronym expansion. Specifically, we start with in-expansion (acronym expansion extraction within documents) in Section 2.1. Then, we move on to out-expansion, which expands acronyms based on other documents and subject matter, in Section 2.2. Finally, we describe the few known end-to-end acronym expander systems in Section 2.3.

2.1 In-expansion

Pustejovsky et al. [59] present a technique that uses a robust parsing of the input text in order to reduce the context within which to search for a candidate expansion. Schwartz and Hearst [63] describe a technique that considers two possible placements of expansions and acronyms in text (before or after), and chooses the correct expansion by matching acronym characters with potential expansion characters.

The MadDog [74] in-expander introduces minimal variations of the Schwartz and Hearst technique [63] which refine the candidate expansions using a sequence of rules where each one yields more precise expansions than the previous one. Nabeesath and Nazeer [62] suggest new pattern heuristics as well as space reduction heuristics. The technique of Azimi et al. [5] uses the same patterns as Schwartz and Hearst [63] but relaxes the heuristics for acronym and expansion extraction: an acronym simply needs to be a token composed of capital letters of some length n and an expansion should be composed of n tokens.

The technique proposed by Yarygina and Vassilieva [82] makes use of user feedback and two decision tree classifiers in order to filter candidate acronym-expansion pairs. Glass et al. [24] proposed a technique that focuses on several languages other than English, and scores candidate pairs by using word embeddings in order to measure the similarity between candidate acronyms and expansions.

In the techniques of Liu et al. [45] and Veyseh et al. [76], the task of finding expansions for an acronym is formalized as a sequence labeling problem solved by Conditional Random Fields (CRFs) [40] based techniques. SciDr [68] in-expander and Zhu et al. [84] also interpret acronym and expansion extraction as a sequence labeling task and make use of pre-trained BERT-based models coupled with ensemble techniques to achieve higher model performance than previous techniques. SciBERT is a language model based on Transformers and pre-trained on research papers from Semantic Scholar¹. SciBERT is fine-tuned in SciDr [68] with training data for the sequence labeling task. The SciDr [68] in-expander uses an ensemble (blending) [67]. It splits the training data into train and validation sets. Five different SciBERT models (e.g., number of epochs and learning rate values) are constructed based on the training set. Predictions on the validation set are then stored. The expansions of the SciBERT models and the rule-based baseline technique of the SDU@AAAI competition² based on Schwartz and Hearst [63], and additional syntactic features extracted from the word-to-tag mapping are used to train five Conditional Random Fields (CRFs) [40] in a 5-fold cross-validation setting. The ensemble technique for these CRF models is based on hard voting.

The technique of Chopard and Spasić [14] also makes use of word embeddings and calculates the *Word Mover's Distance* [39] in order to select the correct expansion from the candidate expansions of an acronym. Jacobs et al. [31] propose a technique that focuses on Modern Hebrew and makes use of a Support Vector Machine (SVM) to select the correct expansion from several candidate expansions for an acronym. Similarly, to select the correct expansion for biomedical documents, in the technique of Kuo et al. [38], an SVM is implemented as well as Logistic Regression and Naïve Bayes models.

Another line of work extracts acronyms not from text but from Web Data like query click logs [32, 72].

Similar tasks are addressed by the fields of Named Entity Recognition and Coreference Resolution. Named Entity Recognition [79] finds entities mentioned in texts and labels them with high level categories like *person* and *organization*; or, for special applications, as molecular biology entities covered in BioNLP tasks [18, 25] like *cells* and *proteins*. Coreference Resolution [52] tries to match expressions like *I*, *my*, or *she* to entities (e.g., names) in text.

2.2 Out-expansion

Classic Context Vector [2, 42, 58] is a typical baseline for out-expansion. It represents the context of an acronym/expansion x by the frequencies of the words in all documents containing x . Li et al. [42] propose two techniques based on word embeddings

from Word2Vec [46] to address the out-expansion problem. Their best technique, called Surrounding Based Embedding, combines the Word2Vec embeddings of the words surrounding the acronym or the expansion. Similarly to Surrounding Based Embedding, Ciosici et al. [16] propose Unsupervised Acronym Disambiguation that replaces each expansion occurrence in a collection of text documents by a normalized token and retrains the Word2Vec google news model [46] on that collection. The resulting model produces an embedding for each normalized token, i.e., an expansion embedding.

Thakker et al. [73] creates document vector embeddings, using Doc2Vec, for each document. For each set of documents D containing an expansion for an acronym A , the system trains a Doc2Vec model on D which is used to infer the embedding for an input document i containing an undefined acronym A .

Charbonnier and Wartena [12] proposed an out-expansion technique based on Word2Vec embeddings weighted by Term Frequency-Inverse Document Frequency scores to find out-expansions for acronyms in scientific document captions.

Recently, works in acronym disambiguation make use of neural networks: MadDog [74] proposes a sequential model to encode context in sentences followed by a feedforward network to classify the input sentence with an expansion. Competitors of the SDU@AAAI competition [75] mainly use pre-trained language models based on Transformer neural networks like BERT [20] and SciBERT [7]. SciDr [68] formulates the out-expansion problem as a substring prediction task. Given a list of expansions concatenated with a sentence as input, it uses the pre-trained language model SciBERT [7] and retrains that model in 5 cross-validations of the sentences dataset to predict the substring, i.e., start and end word indices corresponding to the predicted expansion. The authors also assemble additional SciBERT models trained on external data. The external data considered is constituted by Wikipedia pages that contain an expansion found in the training data.

A related line of work explored the expansion of acronyms in enterprise texts [23, 43]. For instance, in Li et al. [43], enterprise textual documents as well as Wikipedia documents are used as training data. Other works explored acronym out-expansion in biomedical domains [44, 47, 48, 53, 59, 71, 77, 78, 83]. In our work, we explore the general acronym expansion problem where the input document domain or source is not previously known.

Entity Disambiguation (ED) (often referred to as Entity Linking) is the task that links an entity found in text by Named Entity Recognition (NER) to a knowledge base, usually Wikipedia pages [50, 65, 66]. This field is analogous to out-expansion because an expansion can be seen as (and in some cases is) a Wikipedia page title. Several techniques have been proposed to address this task. The survey [65] identifies the work of [81] that is part of the LUKE project³ as the best or one of the best on several datasets, some based on Wikipedia. LUKE (Language Understanding with Knowledge-based Embeddings) [80] is a pre-trained language model that learns to predict masked words and entities. LUKE also employs a global model that, given a set of entities in a document, assigns a ranking among these entities based on confidence.

Moreover, Entity Disambiguation works have explored Natural Language Techniques that we also used in order to represent

¹<https://www.semanticscholar.org/>

²https://github.com/amirveyseh/AAAI-21-SDU-shared-task-1-AI/blob/master/code/character_match.py

³<https://github.com/studio-ousia/luke>

documents like Term Frequency–Inverse Document Frequency (TF-IDF) [35] in [13], Latent Dirichlet Allocation (LDA) [9] in [57], and Doc2Vec [41] in [64, 85].

At BioNLP Open Shared Tasks 2019, Bacteria Biotope [10] considers the goal of linking microbial taxa, habitats, and phenotype to biological knowledge bases. To enrich the input, the authors provided the in-expansions for the acronyms found in their dataset using Ab3p [69]. The winner [36] matched the Word2Vec embeddings of entities in the text with the concepts in the knowledge base. However, an acronym as an entity mention would have the same Word2Vec embeddings regardless of the document.

Less directly related, but insightful, is the literature on Word Sense Disambiguation (WSD) [49, 51] because that work also must make use of the context around a token (in our case, an acronym; in the word sense literature, a word). Raganato et al. [60] proposed a benchmark for word sense disambiguation.

2.3 End-to-end Acronym Expanders

To our knowledge, systems that expand acronyms use a pre-defined dictionary of acronym-expansions [1, 26] as opposed to trying to discover the proper expansion based on context.

Only two end-to-end systems use context for out-expansion. First, Ciosici and Assent [15] propose an end-to-end abbreviation/acronym expansion system architecture that performs out-expansion. Unfortunately, their demo paper provides few technical details and their code is proprietary.

Most recently, Veysseh et al. [74] proposes an end-to-end acronym expansion system, called MadDog, which contains a rule-based expander technique that improves on [63] and an out-expander based on neural networks: a sequential model to encode context followed by a feedforward network to classify the input with an expansion. They also trained their models on a large corpus of sentences.

Neither of these two systems provides a framework with easy plug-in for different in and out-expansions techniques nor uses other data sources. Moreover, neither was evaluated on an end-to-end acronym expander benchmark.

3 AcX: an End-to-end Acronym eXpander System

The AcX system (see Figure 1) consists of:

- (i) A *Database Creation* process which generates an *Expansion Database*⁴ that contains documents, acronyms and their corresponding in-expansions. The Expansion Database also associates each <acronym, in-expansion> pair with a *representation* of the document where that acronym and in-expansion were found. The representation characterizes the content of the document in some summary form. To support other domains and languages, we pass documents in the desired domains/languages to the Database Creation process. In Section A.2 of Annex A, we present additional detail about data structures used in AcX database.
- (ii) The *Acronym Expander Server* that accepts one document at a

⁴When benchmarking, the expansion database will provide us with both a training set and a test set.

time from a user and outputs a list of acronyms found in the input document and the corresponding expansions found by the system (whether as in-expansions or as out-expansions).

For each document with in-expansions, the *Database Creation* process runs the following pipeline:

- (1) an *Acronym and In-Expansion Extractor* obtains the <acronym, expansion> pairs from the document using only within-document evidence.
- (2) a *Representator* (there are many possible representators e.g., Latent Dirichlet Allocation that output topics) maps the document to a document representation that holds document contextual information.
- (3) the *Expansion Database* stores the in-expansions, acronyms, and document representations on disk, currently SQLite [28].

Given a new input document d supplied by a user, the *Acronym Expander Server* executes the following pipeline:

- (1) applies the *Acronym and In-Expansion Extractor* used to build the Expansion Database to extract all the acronyms having expansions in the input document d .
- (2) when d contains links to web pages then those page texts are extracted and inspected to find the expansion for acronyms whose expansions are not found in d (*Link Follower*).
- (3) utilizes the same *Representator* (say, topics from Latent Dirichlet Allocation) used to characterize each document in the Expansion Database to map d to a document representation.
- (4) for each acronym A having no in-expansion in d , the server runs the *Out-Expansion Predictor* to choose a context-appropriate out-expansion. Formally, an expansion E is selected for an acronym A in d if the representations of the documents $doc(A, E)$ with expansion E share more characteristics with the representation of d by some criteria (e.g., closest cosine similarities or labeled by some machine learning classifier for A) than the documents in $doc(A, E')$ for every alternative expansion E' . Thus, for example, if the context of d is publishing, then "PDF" should likely expand to "Portable Document Format" but if the context of d is probability or statistics, then "PDF" should expand to "probability distribution function."

For a language other than English, the in- and out-expansion techniques should be tuned to the new language. They may benefit from changing preprocessing steps such as tokenization for the new language or from adopting a language model trained on the new language or even adopting a multilanguage model.

3.1 Acronym and In-Expansion Extraction

Acronym and in-expansion extraction can use rule-based or machine learning technique. Currently, we have integrated the in-expander implementations of Schwartz and Hearst [63], MadDog in-expander [74], SciBERT [68] and SciDr in-expander [68]. In our rule-based implementations (i.e. Schwartz and Hearst [63] and MadDog [74]), we used roughly the following three-step process as described in [54]:

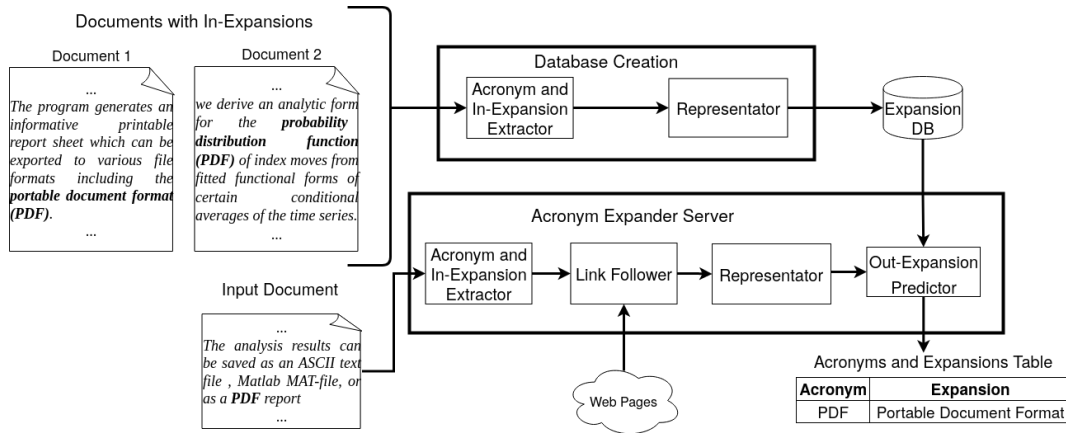


Figure 1: Acronym eXpander (AcX) system. The top stream denotes the creation of the Expansion Database that associates each \langle acronym, in-expansion \rangle pair with some representation of the document(s) where that pair was found. The bottom stream shows the processing of an input document d by combining acronym in-expansion when possible and a representation of d . For an acronym A with no expansion in d , the representation of d is compared with the representations in the Expansion Database of documents containing A to find the context-appropriate expansion.

- (1) *Acronym extraction*: identifies acronyms in a document, e.g., PDF in Figure 1. We modified Schwartz and Hearst [63] to find candidate acronyms even when there is no expansion found in a given document. The technique excludes tokens in which all alphabetic characters except the first character are lower case. We also reject acronyms of two characters where the first is a letter and the second is a dot "." to avoid person names.
- (2) *Candidate expansion extraction*: builds candidate pairs of acronyms and possible in-expansions \langle acronym, expansion \rangle from information in the document, e.g., \langle PDF, formats including the portable document format \rangle from Document 1 in Figure 1.
- (3) *Candidate refinement*: evaluates each candidate pair using a variety of heuristics (e.g., find the shortest expansion that matches the acronym) to obtain a final in-expansion for each acronym that has at least one candidate in-expansion within the document, e.g., portable document format from \langle PDF, formats including the portable document format \rangle .

For the in-expanders of SciBERT and SciDr, the extraction of acronyms and expansions is formalized as a sequence tagging problem where each token can have one of three tags: (i) a token in an acronym (e.g., CD in CD-ROM), (ii) a tokenword in an expansion, or (iii) other token. For example, from Document 1 in Figure 1, PDF would be tagged as an acronym token, each token portable, document, and format would be tagged as a token in an expansion. The remaining tokens in Document 1 would have the "other token" tag. AcX appropriately tags the acronym-expansion pairs and other tokens in the training data, then builds a machine learning model on the tagged data. The output of such machine learning models is then converted to acronym-expansion pairs by matching the acronym characters against expansions.

Our system supports ensemble in-expansion through SciDr. That ensemble technique can be easily extended to include additional in-expansion techniques.

3.1.1 Link Follower For an input document d containing hyperlinks, the *Link Follower* component follows those links to try to find the expansions from documents pointed to by d . This module is used after executing the acronym and in-expansion extraction algorithm for the acronyms with no expansion found in text. For each A acronym having no in-expansion, the Link Follower is constituted by the following steps, until an expansion is found:

- (1) searches in the *title* attribute of the hyperlink an HTML tag for the expansion. For example, given the following tag $\langle a \ href="/wiki/Leadership_in_Energy_and_Environmental_Design" \ title="Leadership \ in \ Energy \ and \ Environmental \ Design" \rangle LEED \langle /a \rangle$ the expansion "Leadership in Energy and Environmental Design" would be extracted for acronym LEED.
- (2) resolves relative HTML URL links to global ones. For example, if no expansion were found in the example above with a local URL, the resolved URL would be https://en.wikipedia.org/wiki/Leadership_in_Energy_and_Environmental_Design.
- (3) executes the same in-expander technique used by the acronym and in-expansion extractor.

Given a title, the acronym is placed inside parenthesis and appended to the end. Then, the in-expander technique is executed on the title in order to extract a possible expansion.

3.2 Representator

Representors in the AcX system summarize documents in order to capture knowledge about their semantics. Although AcX supports sentence-level out-expansion techniques, using the whole document is more effective than using just parts of the text because the whole document captures the overall context better.

Some representators assign a set of topic terms description to a document. If two documents have many topic representative terms in common, then they are considered to be semantically related.

Other representators use embeddings [41] to characterize a document. An *embedding* is a vector of real numbers in a high dimensional space. Embedding techniques map an object encoded in a one-hot representation, a very sparse and high dimensional vector of binary values, into a very dense and lower dimensional vector of real values (i.e., embedding). A small distance between embedding vectors suggests document similarity.

AcX encloses several techniques that can semantically represent an entire set of documents that contain the same expansion for a given acronym. Specifically, let $docs(A, E)$ denote the set of full document texts in which a given acronym A is defined by a single expansion E (e.g., all documents in which acronym PDF is explicitly expanded as portable document format):

Here are some representations of such a collection of documents:

- *Classic Context Vector (CCV)* [2], represents an expansion E by the set of words in $docs(A, E)$ along with their counts.
- *Document Context Vector (DCV)* (our variation of context vector), builds on context vector, however it represents each document $d \in docs(A, E)$ individually by the set of word occurrences in d . For example, the word occurrences corresponding to Document 2 in Figure 1 would contain among others the values {of: 3}, {the: 2}, {derive: 1}, {analytic: 1}, {form: 1}.
- *Term Frequency–Inverse Document Frequency (TF-IDF)* [35], weights each term t in each document $d \in docs(A, E)$ highly if it is found frequently in d and infrequently in the entire document corpus, thus permitting the characterization of each document by its highly weighted terms. For example, the TF-IDF score for the word the in Document 2 in Figure 1 is $\frac{2}{27} \cdot \log(\frac{2}{2}) = 0$ because this word appears in both documents.
- *Latent Dirichlet Allocation (LDA)* [9] assigns topics to documents using a Dirichlet probabilistic model. For example, Document 2 in Figure 1 could be represented by the following topics: topic1={{analytics: 0.7}, {series: 0.3}} and topic2={{functional: 0.8}, {form: 0.2}}.
- *Doc2Vec* [41] is a document embedding technique based on Word2Vec [46] which assigns vectors to words in such a way that words that appear in the same context have a high cosine similarity. For example, the words functional and conditional would be assigned similar vectors. Thus, using the principles of Word2Vec, Doc2Vec assigns vectors to entire documents. For example, documents 1 and 2 in Figure 1 would be assigned mutually distant vectors.
- *Sentence Bidirectional Encoder Representations from Transformers (SBERT)* [61] constructs sentence embeddings that can be compared to determine sentence similarity. AcX splits the input document text to fit into the SBERT input limit (e.g., 384 tokens), and then we average the resulting embedding vectors to get a document representation.

3.3 Out-Expansion Predictor

To choose an out-expansion for an acronym A in an input document d having no expansion for A , the Out-Expansion Predictor component considers each candidate out-expansion E for A and compares d to some representation of $d' \in docs(A, E)$.

In the case of Classic Context Vector (CCV), we compare d with the vector representation of $docs(A, E)$. For the remaining techniques, we compare d with each document representation of $d \in docs(A, E)$.

Using cosine similarity, the Out-Expansion Predictor will choose an out-expansion E over a different expansion E' if any document $d \in docs(A, E)$ is more similar to d than all $d''' \in docs(A, E')$.

A classical similarity technique is cosine similarity, but the AcX system also supports classification-based approaches that work as follows. Consider all the documents, denoted $alldocs(A)$ containing in-expansions of acronym A . Some documents in $alldocs(A)$ have an in-expansion of $E1$ for A , some have $E2$ for A and so on. Given the representations of documents in $alldocs(A)$ as features and the expansions ($E1, E2$, etc) as labels, the out-expansion problem becomes a machine learning classification problem. When a new document d is given to AcX, the representation of d is passed as input (i.e., features) to the classifier which labels d with an expansion.

The classifiers we support so far are:

- *Support Vector Machines (SVMs)* [19] fit a hyper-plane that optimally separates binary labeled data in the feature space. Non-binary classification is performed by a "one-vs-all" technique where a binary SVM classifier predicts with a certain probability if an input document belongs to a particular class (where each class corresponds to a particular expansion). The class (and therefore expansion) with the highest probability is selected. We used the LibLinear [22] implementation included in scikit-learn toolkit [56].
- *Logistic Regression (LR)* [33] fits a logistic function to classify binary classes (again a class corresponds to an expansion). Non binary classification is again performed by a "one-vs-all" technique. We used the LibLinear [22] implementation included in scikit-learn toolkit [56].
- *Random Forests (RF)* [11] fit a particular number of decision trees (default 100) trained on randomly selected samples. There will be one random forest per acronym A . The representation of a document having no in-expansion for A will be input to the random forest. Each tree will predict one expansion with some probability. The random forest selects the class whose average probability is the highest. We used the scikit-learn [56] implementation.

In addition to these classifiers, for evaluation purposes or for anyone who wants to try other techniques, AcX supports the following additional techniques from related work: Surrounding Based Embedding (**SBE**) [42], Thakker et al. [73], Unsupervised Abbreviation Disambiguation (**UAD**) [16], the SciDr out-expander (**SciDr-out**) [68], the MadDog out-expander (**MadDog-out**) [74], and **LUKE** [81], a state-of-the-art technique for Entity Disambiguation. For UAD, SciDr-out and MadDog-out, AcX performs sentence segmentation and, given the results from each sentence, decides which

expansion to assign to the text. For UAD, we select the most frequent predicted expansion among the sentences in the document.

We have extended SciDr-out to consider all the sentences containing the acronym A instead of just one sentence as in SciDr-out’s original implementation. SciDr-out associates an acronym with its possible expansions concatenated together. The system then finds the substring of that concatenated string with the highest probability and outputs that as the expansion. For example, the concatenated expansion of "PDF" might be "probability density function portable document format". Depending on the contents of some input document d containing "PDF", SciDr-out will choose some substring of that concatenated expansion.

We have extended MadDog-out to enable it to train in new documents, instead of using only their original machine learning models. MadDog-out processes the last sentence of any document containing acronym A to determine the most likely expansion.

For LUKE, we had to modify the internals to work with acronyms and expansions. We use their pre-trained model and perform fine-tuning in our training data using the procedure described by the authors in [81], except that we allow the entity embeddings (now expansion embeddings) to be updated during training. This modification allows the generation of embeddings for expansions out of the original model vocabulary.

4 In-expansion Benchmark, Evaluation and Results

We describe our benchmark of in-expansion techniques in Section 4.1 and evaluate state-of-the-art techniques on this benchmark in Section 4.2.

4.1 A Benchmark of In-expansion Techniques

This section describes the benchmark we developed to evaluate in-expansion techniques. Section 4.1.1 details the datasets used in this benchmark. Section 4.1.2 lists the in-expansion techniques that we implemented for this benchmark. Section 4.1.3 defines the metrics that we used to evaluate the in-expansion extraction techniques.

4.1.1 Datasets The datasets included in this in-expansion benchmark are:

Medstract: This dataset is composed of 199 randomly selected MEDLINE⁵ abstracts from the results of a query on the term "gene". The abstracts were manually annotated and then the annotations were corrected and improved by Schwartz and Hearst [63], Ao and Takagi [3], Pustejovsky et al. [59], Yarygina and Vassilieva [82] and Doğan et al. [29]. We use the last revised version of Doğan et al. [29] that contains 159 acronym-expansion pairs.

Schwartz and Hearst: This dataset consists of 1,000 randomly selected MEDLINE abstracts from the results of a query on the term "yeast". The abstracts were manually annotated by Schwartz and Hearst [63] and revised by Doğan et al. [29]. The revised version that we use contains 979 acronym-expansion pairs.

⁵<https://www.nlm.nih.gov/bsd/medline.html>

BIOADI: This dataset contains 1,201 abstracts from the BioCreative II gene normalization dataset. The dataset was original annotated by Kuo et al. [38] and revised by Doğan et al. [29]. It contains 1,720 acronym-expansion pairs.

Ab3P: This dataset results from the random selection of MEDLINE 1,250 abstracts. The dataset was manually annotated by Sohn et al. [69]. We use the revised version of Doğan et al. [29] that contains 1 223 acronym-expansion pairs.

SciAI: This dataset results from processing 6,786 English arXiv⁶ papers. Those papers were split into sentences and sent to Amazon Mechanical Turk (MTurk) to be annotated by humans, resulting in 9,775 acronym-expansion pairs. This dataset was annotated for both acronyms and acronym-expansion pairs. The final dataset has 17,506 sentences, where 1% do not contain acronyms and 24% do not contain expansions. We use the SDU@AAAI competition [75] version⁷ that was initially proposed by Veysseh et al. [76].

End-to-end: We developed a dataset that consists of 163 English Wikipedia documents randomly selected from the *Computing* category⁸ in Wikipedia. It contains 1,139 acronym-expansion pairs. Although intended to evaluate an end-to-end acronym-expander system, for this in-expansion benchmark in particular, we consider only the acronym-expansion pairs with expansion in text. (Later, in Section 6.1, we use the whole set of acronym-expansions pairs to evaluate end-to-end systems.) Each document was annotated by two computer science students who volunteered for the task. Each student annotated at least two documents. During the annotation process, each student identified each acronym in the document and mapped it to an expansion. Each acronym-expansion pair was labeled by the annotators, indicating whether the expansion was present in text. Any conflict between annotators was manually resolved by the authors. The Inter-Annotator Agreement (IAA) among each annotators (excluding the third annotator, the reviewer) using Krippendorff’s alpha [37] with the MASI distance metric [55] is 0.68 for in-expansion pairs and 0.33 for out-expansion pairs. In a hypothetical scenario, if both annotators had given the same acronym-expansions, then the score would be 1. In this case, the human annotators disagree on out-expansions more often than on in-expansions. This is unsurprising because out-expansion requires consulting additional text sources other than the document at hand, while for in-expansion the text provided is enough. In Section A.1 of Annex A, we present additional detail about the process used to create this dataset.

4.1.2 In-expansion techniques This benchmark includes the following in-expansion techniques (that are supported by our AcX system described in Section 3):

Rule-based: Schwartz and Hearst (SH) [63] technique and the MadDog [74] in-expansion (**MadDog-in**) technique which builds on the Schwartz and Hearst algorithm.

⁶<https://arxiv.org/>

⁷<https://github.com/amirveysseh/AAAI-21-SDU-shared-task-1-AI>

⁸<https://en.wikipedia.org/wiki/Category:Computing>

Machine Learning: SciBERT based technique used in [68] and the SciDr [68] in-expansion (**SciDr-in**) technique which ensembles SciBERT models and a rule-based technique based on SH with Conditional Random Fields. Moreover, we consider models used by these machine learning techniques that are trained *with external data* besides the individual training sets of each dataset. The external data is composed of Medstract, Schwartz and Hearst, BIOADI, and Ab3P train sets if the test set is biomedical. For SciAI and End-to-end test sets, the external data consists of all train sets (i.e., biomedical datasets, SciAI, and End-to-end).

4.1.3 Performance metrics Our benchmark uses the following metrics. The metrics apply to acronyms alone as well as to acronym-expansion pairs. The acronyms can be either in singular or plural form to be considered equal, and the expansions are equal if their lower case versions without dashes have an edit distance less than 3 or if the first 4 characters of each word are equal. If the same acronym or pair appears several times in the same document, it is counted only once:

Acronym Pair Precision: the number of correctly extracted acronym pairs divided by the number of acronym pairs extracted by that technique over all documents. It will be calculated as:

$$\frac{\# \text{ of correctly extracted acronym pairs}}{\# \text{ of extracted acronym pairs}}$$

Acronym Pair Recall: the number of correctly extracted acronym pairs divided by the number of distinct acronym pairs present over all documents. It will be calculated as:

$$\frac{\# \text{ of correctly extracted acronym pairs}}{\# \text{ acronym pairs in text}}$$

Acronym Pair F1-measure: the harmonic mean of the *precision* and *recall* of the system. It will be calculated as:

$$2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Acronym Pair Cohen's Kappa: measures the agreement between the acronym pairs extracted by a technique and the actual acronym pairs present in text for a given number of documents in a dataset. The value of *kappa* [17] is calculated as $\frac{p_o - p_e}{1 - p_e}$, where p_o is the probability of agreement between an acronym-expansion extraction technique and the actual dataset and p_e is the probability of random agreement. This measure is particularly useful for unbalanced test sets, because it gives greater weight to infrequent acronym-expansions (for which the denominator will be smaller). To avoid treating acronyms and expansions with small differences as entirely different labels (i.e. full disagreement), we calculate kappa with respect to the *Jaccard* [30] distance.

Training time: CPU or GPU time in seconds to train the machine-learning models that are used by the in-expansion technique.

Execution time: CPU or GPU time in seconds that the in-expansion technique takes to extract acronym-expansion pairs from a document in the dataset.

By default (i.e., when not mentioned otherwise), acronym pair *Precision*, *Recall*, and *F1-Measure* are calculate on the *Document-Level* where every unique acronym-expansion pair (or acronym) extracted from each document is counted. Thus, an acronym-expansion pair that appears many times across documents in the corpus will

be weighted more than one that appears less often. This is in line with the most recent efforts in in-expansion which aim to efficiently extract acronym-expansions from each individual document independently.

Additionally, we calculate the *Precision*, *Recall*, and *F1-Measure* alternative level, *Glossary*. In the *Glossary-level*, an acronym-expansion pair (or just acronym) that appears in the text or is expanded multiple times in different documents in the corpus is counted only once. This follows the first measurements performed in-expansion techniques which aimed to create only a global glossary (dictionary) based on a set of documents.

To better explain the *Document* and *Glossary* levels, let us consider a dataset D that is composed of documents A and B . In document A there are five instances of the acronym-expansion pair x and in document B two, with D having seven instances of x . Furthermore, a system for acronym-expansion extraction is run on dataset D and only extracts the five instances of x in document A . On the *Glossary level*, the number of correctly extracted pairs and the total number of extracted pairs are equal to one. Furthermore, the number of total acronym-expansion pairs present in D is only counted as one, because only unique pairs are considered. This results in a *pair precision* of $\frac{1}{1} = 1$ and a *pair recall* of $\frac{1}{1} = 1$. On the *document level*, the number of correctly extracted pairs and total number of extracted pairs is five. However, the number of total acronym-expansion pairs present in D is counted as seven, due to the five instances of x in document A and two instances in document B . This results in a *pair precision* of $\frac{5}{5} = 1$ and a *pair recall* of $\frac{5}{7} = 0.7$.

4.2 In-expansion Experimental Evaluation

In this section, we evaluate the in-expansion techniques using the benchmark presented in Section 4.1.

Setup. The in-expansion experiments were performed on a machine with an Intel® Core™ i5-4690K CPU with 4 cores, and 16 GB of RAM and an NVIDIA GeForce GTX 1070. Only SciBERT and SciDr-in used the GPU.

Results.

We report Precision, Recall, F1-measure, and Cohen's kappa (K) values for the biomedical datasets (i.e., Medstract, Schwartz and Hearst, BIOADI and Ab3P) for acronym extraction in Table 1 and for acronym-expansion pair extraction in Table 2. We report Precision, Recall, F1-measure, and Cohen's kappa (K) values for both acronym and pair for SciAI dataset in Table 3 and for End-to-end dataset in Table 4. The additional external data used to train SciBERT and SciDr-in for the biomedical application includes the data of all biomedical datasets excluding the test set (30%). For SciAI and End-to-end datasets, the external data used to train SciBERT and SciDr-in includes all documents in the other datasets (i.e., Medstract, Schwartz and Hearst, BIOADI, Ab3p, SciAI, and End-to-end).

For the *biomedical domain*, for all acronym and pair extraction measures (precision, recall, and F1-measure), the best acronym and in-expander technique, on average, is SH.

On *SciAI*, the machine-learning techniques SciDr-in and SciBERT outperform the rule-based techniques, SH and MadDog-in, in terms of recall (**91.78%-94.05%** for acronym and **88.24%-90.64%**

Acronym and In-expansion Technique	Acronym																			
	Medstract				SH				BIOADI				Ab3P				Averages			
	P	R	F1	K	P	R	F1	K	P	R	F1	K	P	R	F1	K	P	R	F1	K
SH	100.00%	89.13%	94.25%	0.97	99.56%	81.13%	89.50%	0.91	99.06%	79.58%	88.25%	0.88	98.62%	77.66%	86.89%	0.89	99.31%	81.88%	89.72%	0.91
MadDog	100.00%	63.30%	77.33%	0.77	96.64%	51.80%	67.45%	0.77	97.99%	55.19%	70.61%	0.74	99.16%	64.31%	78.02%	0.85	98.45%	58.65%	73.35%	0.78
SciBERT	81.25%	56.52%	66.67%	0.68	85.56%	70.50%	77.47%	0.77	88.06%	73.91%	80.37%	0.76	85.99%	71.93%	78.34%	0.78	85.22%	68.22%	75.71%	0.75
SciBERT with External Data	86.84%	71.74%	78.57%	0.76	88.16%	77.70%	82.60%	0.82	90.61%	78.45%	84.09%	0.82	87.46%	76.02%	81.34%	0.81	88.27%	75.98%	81.65%	0.80
SciDr	93.33%	60.87%	73.68%	0.73	87.56%	60.79%	71.76%	0.74	92.11%	64.08%	75.58%	0.75	89.25%	67.85%	77.09%	0.79	90.56%	63.40%	74.53%	0.75
SciDr-in with External Data	92.11%	76.09%	83.33%	0.83	93.13%	78.05%	84.93%	0.86	93.31%	79.21%	85.69%	0.84	89.08%	71.11%	79.09%	0.81	91.91%	76.12%	83.26%	0.84

Table 1: In-expansion techniques Precision, Recall, F1-measures and Cohen’s kappas (K) for acronym for each biomedical dataset (Medstract, SH, BIOADI, and Ab3p) and the averages.

Acronym and In-expansion Technique	Pair																			
	Medstract				SH				BIOADI				Ab3P				Average			
	P	R	F1	K	P	R	F1	K	P	R	F1	K	P	R	F1	K	P	R	F1	K
SH	100.00%	89.13%	94.25%	0.97	96.03%	78.42%	86.34%	0.89	94.11%	75.61%	83.86%	0.83	95.16%	74.93%	83.84%	0.86	96.33%	79.52%	87.07%	0.89
MadDog	93.10%	58.69%	72.00%	0.35	93.29%	50.00%	65.11%	0.47	87.58%	49.33%	63.12%	0.48	95.37%	61.85%	75.04%	0.64	92.34%	54.97%	68.82%	0.49
SciBERT	65.62%	45.65%	53.84%	0.27	72.37%	59.35%	65.21%	0.41	67.79%	56.90%	61.87%	0.36	74.27%	62.13%	67.66%	0.53	70.01%	56.01%	62.15%	0.39
SciBERT with External Data	76.32%	63.04%	69.05%	0.33	76.32%	67.26%	71.51%	0.42	74.45%	64.46%	69.10%	0.44	76.80%	66.76%	71.43%	0.54	75.97%	65.38%	70.27%	0.43
SciDr	80.00%	52.17%	63.16%	0.38	74.61%	51.79%	61.14%	0.4	76.90%	53.50%	63.10%	0.41	81.00%	61.58%	69.97%	0.54	78.13%	54.76%	64.34%	0.43
SciDr-in with External Data	92.11%	76.09%	83.33%	0.39	83.69%	70.14%	76.32%	0.49	86.19%	73.16%	79.14%	0.5	80.20%	64.03%	71.21%	0.52	85.55%	70.86%	77.50%	0.48

Table 2: In-expansion techniques Precision, Recall, F1-measures and Cohen’s kappas (K) for pair for each biomedical dataset (Medstract, SH, BIOADI, and Ab3p) and the averages.

Acronym and In-expansion Technique	SciAI							
	Acronym				Pair			
	P	R	F1	K	P	R	F1	K
SH	96.02%	82.36%	88.67%	0.97	92.85%	79.64%	85.74%	0.88
MadDog-in	98.63%	86.72%	92.30%	0.98	96.91%	85.21%	90.68%	0.96
SciBERT	95.69%	94.05%	94.86%	0.97	92.21%	90.64%	91.42%	0.94
SciBERT with External data	96.18%	94.05%	94.90%	0.97	92.50%	90.45%	91.46%	0.93
SciDr-in	97.47%	92.47%	95.11%	0.98	94.47%	89.63%	91.98%	0.95
SciDr-in with External data	97.58%	91.78%	94.59%	0.98	93.81%	88.24%	90.94%	0.94

Table 3: In-expansion techniques Precision, Recall, and F1-measures and Cohen’s kappas (K) for acronym and pair extraction and for the SciAI dataset.

Acronym and In-expansion Technique	User Generated							
	Acronym				Pair			
	P	R	F1	K	P	R	F1	K
SH	91.00%	70.54%	79.47%	0.72	86.00%	66.67%	75.10%	0.14
MadDog-in	92.78%	69.76%	79.64%	0.71	88.65%	66.67%	76.10%	0.11
SciBERT	65.62%	48.83%	55.99%	0.62	58.34%	43.41%	49.77%	0.08
SciBERT with External data	49.67%	58.91%	53.90%	0.62	45.09%	53.48%	48.93%	0.08
SciDr-in	77.08%	57.36%	65.77%	0.66	68.75%	51.16%	58.66%	0.09
SciDr-in with External data	86.36%	58.91%	70.04%	0.69	81.81%	55.81%	66.35%	0.11

Table 4: In-expansion techniques Precision, Recall, F1-measures and Cohen’s kappas (K) for acronym and pair extraction and for the User Generated dataset.

for pair) and F1-measure (94.59%-94.05% for acronym and 90.94%-91.98% for pair) for both acronym and pair extraction. Furthermore, MadDog-in achieves the best precisions (98.63% for acronym and 96.91% for pair) followed by SciDr-in (97.47%-97.58% for acronym and 93.81%-94.47% for pair).

On the *End-to-end dataset*, MadDog-in achieves the best overall precision (92.78% and 88.65%) and F1-measure (79.64% and 76.10%) for acronym and pair extraction, while SH surpasses MadDog-in in terms of acronym recall (70.54%) and matches for pair recall (66.67%). Furthermore, among the machine-learning techniques on the *End-to-end dataset*, SciDr-in surpasses SciBERT on precision (77.08% and 68.75%), recall (57.36% and 51.16%), and F1-measure (57.36% and 58.66%) for both acronym and pair extraction. Increasing the training dataset with External Data for SciBERT yielded an increase in recalls (58.91% and 53.48%) but decreased precisions (49.67% and 45.09%) and F1-measures (53.90% and 48.93%) for both acronym and pair extraction, while for SciDr-in, we observe a general increase in performance.

Cohen’s Kappa analyzes. In general, the Cohen’s Kappa (K) best technique for each dataset is consistent with the F1-measure in the biomedical datasets, i.e., SH technique is clearly the best. On the *SciAI dataset*, MadDog-in scores the best K for acronyms and pairs with 0.98 and 0.96, respectively.

On the *End-to-end dataset*, all techniques achieve values of Cohen’s Kappa higher or equal to 0.62 for acronym extraction and lower or equal to 0.14 for pair extraction. The reason is that the end-to-end data contains many distinct expansions with different frequencies, leading to lower scores.

Glossary-level differences. In Annex B, we report the results obtained using the Glossary-level metrics for each dataset. Regarding the glossary-level results, they are globally very similar to the document-level ones in terms of the differences among in-expansion techniques. The exceptions in the best methods for each metric are small differences: On the BIOADI dataset, the best score for acronym extraction recall is obtained with SciDr-in trained in all biomedical datasets with 79.71%, SH (which was the best on document-level) scores 78.26%. On the SciAI dataset, acronym recall is better for SciBERT trained on SciAI training set with 93.58%, SciBERT

Acronym and In-expansion Technique	Train models execution times (s)							
	Train Dataset						with External Data	
	Medstract	SH	BIOADI	Ab3P	SciAI	User-Generated	All Biomedical	All
SciBERT	108	745	806	831	1 701	421	2 691	5 122
SciDr-in	1 226	7 255	8 060	8 738	52 257	2 752	28 612	99 025

Table 5: In-expansion train models execution times for each dataset and when trained with External Data.

Average execution times per document (s)							
Acronym and In-expansion Technique	Medstract	SH	BIOADI	Ab3P	SciAI	User Generated	Average
SH	0.00	0.02	0.01	0.01	0.05	0.00	0.02
MadDog-in	0.60	4.33	7.23	5.04	15.75	4.53	6.25
SciBERT	10.35	70.29	115.23	80.81	686.44	35.97	166.52
SciBERT with External Data	16.80	73.45	122.73	76.80	645.70	42.84	163.05
SciDr-in	165.69	1 039.45	1 831.67	1 225.03	11 105.72	470.88	2 639.74
SciDr-in with External Data	172.81	1 073.44	1 836.02	1 198.92	11 234.76	486.81	2 667.13

Table 6: In-expansion average execution times per document for each dataset.

trained on all datasets (which was the best on document-level) scored **93.49%**. On the User Generated dataset, the best technique for acronym and the pair extraction F1-measure is SH with **79.49%** and **66.67%** respectively, while Maddog-in (the best on document-level for acronym F1 and one of the best for pair F1) scored **79.04%** for acronym F1, and **65.83%** for pair F1.

Interpretation: In this in-expansion benchmark, rule-based techniques SH and MadDog-in generally perform best for all datasets. The one exception is on the SciAI dataset where machine learning techniques from SciDr-in and SciBERT work better.

Rule-based systems work well for in-expansion, because acronyms follow human-understood rules, viz. roughly, acronyms should be in upper-case, each letter should represent a word, and the expansion should either precede or follow the first use. So it is natural that a rule-based system would do well. Machine learning work better when given more examples (SciAI dataset), however even ensembled with a rule-based technique (SciDr) the results were generally inferior to using the rule-based technique by itself.

While the expansions found by the rule-based techniques are not a superset of those found by the machine learning techniques, SciDr often fails because it adds extra words to the expansion string. On the other hand, SciDr can find unusual cases where not all acronym chars belong in the expansion, e.g., expansion *PIN-FORMED* of *pin1*.

Execution time analysis. We report in Table 6 the in-expansion execution times per document. We observed from our experiments that the rule-based techniques are much faster than the machine learning techniques. SH is the fastest technique on every single dataset taking less than **0.06** seconds on average to extract acronym-expansion pairs from a document. MadDog-in is the second fastest technique taking up to **16** seconds to process each document. The machine-learning techniques SciDr-in and SciBERT take much more time to extract pairs from the datasets than the rule-based techniques. For instance, on the SciAI dataset, SciBERT takes **687** seconds. Comparing SciDr-in and SciBERT, the execution times per document of SciDr-in are much higher than SciBERT taking **11 106** seconds on the SciAI dataset because it is an ensemble

based technique. We report in Table 5 the train execution times for machine-learning based techniques (SciBERT and SciDr-in) for each dataset and with External Data. Regarding training times, SciBERT takes **1 700** seconds on SciAI training data, **2 691** seconds on the biomedical datasets, and **5 121** seconds to train with all datasets. SciDr-in takes longer for training: **52 257** seconds on SciAI training data, **28 612** seconds on the biomedical datasets, and **99 024** seconds to train with all datasets.

In summary:

- If document processing needs to be fast and there are hardware limitations, the rule-based techniques **SH** and **MadDog-in** are the best.
- If there are no time constraints and a large volume of data from the same domain is available, the machine learning techniques **SciBert** and **SciDr-in** are marginally better.
- If the dataset is sentence-based and a large amount of data from the same domain is available, use **SciDr-in**, though, in the medical domain, **SH** would likely be a strong contender.

5 Out-expansion Benchmark, Evaluation and Results

We describe our benchmark of out-expansion techniques in Section 5.1 and evaluate state-of-the-art techniques on this benchmark in Section 5.2.

5.1 A Benchmark of Out-expansion Techniques

This section presents our benchmark of out-expansion techniques. Section 5.1.1 describes the datasets used in this benchmark. Section 5.1.2 explains the steps used to prepare those datasets. Section 5.1.3 lists the out-expansion techniques included in the benchmark, grouped by type. Finally, Section 5.1.4 describes the metrics to evaluate those out-expansion techniques.

5.1.1 Datasets The datasets included in our out-expansion benchmark are:

Statistics	SciWISE	MSH	CS Wiki	SciAD
# of articles	4 677	12 053	10 220	39 815
Average # of chars per article	1 193	1 552	9 246	187
# of sentences	40 300	112 456	702 387	51 340
# of distinct acronyms	129	67	630	732
# of distinct ambiguous acronyms	100	64	566	682
Average # of distinct expansions per article	1.09	0.99	1.02	1
# of distinct acronym/expansions	272	139	8 617	2173
Average # of expansions per ambiguous acronym	2.43	2.13	15.11	3.11

Table 7: Statistics of the out-expansion datasets. SciAD is the dataset with the largest number of articles, 39k. However, if we compare the number of sentences found in each dataset, CS Wiki has the largest total with 70K, 11k for MSH, 5k for SciAD and 4k for ScienceWISE. The reason is that SciAD contains the smallest articles in terms of characters (mostly just sentences having fewer than 200 characters), while others have 1k or more. CS Wiki has the biggest set of distinct acronym-expansions pairs (8k). In terms of distinct expansions per ambiguous acronym, CS Wiki has around 15, while the remaining dataset have maximum around 3.

MSH dataset [34] contains biomedical document abstracts from the MEDLINE (Medical Literature Analysis and Retrieval System Online) corpus used in Li et al. [42], Prokofyev et al. [58]. This dataset was automatically annotated using citations from MEDLINE and the ambiguous terms with MeSH headings identified in the Metathesaurus⁹. We use the original texts and the revised labels from Li et al. [42];

SciWISE dataset consists on the Physics dataset used in Li et al. [42] and Prokofyev et al. [58] that consists of document abstracts. This dataset was annotated by human experts, and it includes expansions either containing at least 2 words or a single word with at least 14 characters.

CS Wiki (Computer Science Wikipedia) dataset created in Thakker et al. [73] contains documents from different fields that contain acronyms used in computer science. Expansions were extracted by parsing the content of English Wikipedia disambiguation pages of acronyms used in computer science (e.g., [https://en.wikipedia.org/wiki/PDF_\(disambiguation\)](https://en.wikipedia.org/wiki/PDF_(disambiguation))).

SciAD This dataset was prepared for the out-expansion SDU@AAAI-21 competition [75]. It is based on the SciAI in-expansion dataset, described in Section 4.1.1. We use the revised version¹⁰ created by Egan and Bohannon [21] who removed duplicate sentences from the original training and validation sets.

Table 7 presents relevant statistics about each dataset. An ambiguous acronym is an acronym having more than one expansion available in the dataset.

5.1.2 Data Preparation The data preparation steps are roughly the same for each out-expansion technique:

⁹https://www.nlm.nih.gov/research/umls/knowledge_sources/metathesaurus

¹⁰<https://github.com/PrimerAI/sdu-data>

- (1) **Dataset Splitting:** We split each dataset into *train* and *test* sets (respectively 70% and 30% of the documents of the original dataset). We then apply 5-fold cross validation on the train dataset in order to tune the hyperparameters of each out-expansion technique. The hyperparameter-tuned technique is then tested on the yet unseen 30% of the data.
- (2) **Expansion Consolidation:** For the expansions of acronym *A* in each dataset, we apply an approximate duplicate detection process that groups expansion strings that correspond to the same expansion meaning. For example, *portable document format* and *Portable-Documents-Formats* are two distinct strings that refer to the same real expansion. As criteria, we consider two expansions to be equal if their lower case versions without dashes have an edit-distance less than 3 or if the first 4 characters of each word are equal. Equal expansions are consolidated by replacing in text all expansion strings with the same meaning by the most frequent expansion.
- (3) **Expansion Removal:** When testing the accuracy of out-expansion techniques on some document *d*, we associate any acronym *A* in the document with its in-expansion *In(A)*, if present. Then, we replace all occurrences of the in-expansion *In(A)* in text by *A* alone.
- (4) **Tokenization:** We apply the word tokenization from the Natural Language Toolkit (NLTK) [8] to obtain only alphanumeric tokens. Additionally, we remove stop words using NLTK and numeric tokens;
- (5) **Token Normalization:** We transform each token into its stem, e.g., probable, probability, and probabilities all map to *probabl*. We use the Porter Stemmer algorithm from NLTK.

The preparation of the MSH and SciWISE datasets follows the preprocessing reported in Li et al. [42], so we apply all the preparation steps above except token normalization. The five steps are consistent with the pre-processing steps used in Thakker et al. [73] for the CS Wiki dataset. For SciDr-out and MadDog-out, we apply only the first three steps, because these techniques replace the last two steps with steps that depend on the language models of the neural networks they use.

5.1.3 Out-expansion Techniques This benchmark includes the following groups of out-expansion techniques:

Classical Techniques: We use two baselines: **Random** which randomly assigns a possible expansion to an acronym; and **Most Frequent** which always selects the most frequent expansion found in our training data as measured by the number of occurrences in distinct documents. We use the Cosine similarity (**CosSim**) with the Classic Context Vector (**CCV**) [42], Document Context Vector (**DCV**) - variant of Classic for each document, Surrounding Based Embedding (**SBE**) [42], and **Thakker** et al. [73].

Sentence-oriented Techniques: We include related work techniques that expect a sentence as input (instead of a document) and adapt them as described in the AcX overview (Section 3.3). These include Unsupervised Abbreviation Disambiguation (**UAD**) [16], MadDog [74] out-expander (**MadDog-out**), and SciDr [68] out-expander (**SciDr-out**).

We also use **SciDr-out with External Data** consisting of the Wikipedia pages that contain an expansion found in the training data.

Representator Techniques: We include **Cossim** with the document representation techniques described in Section 3.2, that we have adapted from natural language processing: Term Frequency-Inverse Document Frequency (**TF-IDF**), Latent Dirichlet Allocation (**LDA**), **Doc2Vec**, and Sentence Bidirectional Encoder Representations from Transformers (**SBERT**). We used SBERT model *all-mpnet-base-v2*, the top performing model in Sentence Similarity tasks (14 datasets)¹¹. *all-mpnet-base-v2*¹² is based on MPNet model [70] that outperforms BERT and RoBERTA in both quality and speed. *all-mpnet-base-v2* was trained on one billion sentences pairs from a diverse set of data sources.

Classification Techniques: We created a complete new class of out-expansion techniques that use the outputs of a representator as features for a Machine Learning classifier, specifically, Random Forests (**RF**), Logistic Regression (**LR**), and Support Vector Machines (**SVM**). Each acronym has its own classifier trained with the features of the documents that contain an expansion for the acronym (e.g., acronym PDF will have a random forest RandFor(PDF) based on documents that contain an in-expansion for PDF). Based on the features of a target document *d*, the classifier will choose the appropriate expansion as explained in Section 3.3.

Combination of Representator Techniques: The final type of out-expansion techniques that we assembled consists of combining two representators' outputs, namely the Doc2Vec with a Context Vector (either Classic or Document), as input to predictors: **CCV + Doc2Vec** and **DCV + Doc2Vec**. Combinations are constructed by concatenating the outputs together into a single feature vector.

Ensembler Techniques: We support two ensembler techniques: **Hard** voting where each technique votes for its preferred expansion regardless of its confidence; and **Soft** voting that takes the averages of confidences per expansion. The confidences are normalized at the individual technique level in such a way that their sum is 1. For the experiments, we assembled the following 7 out-expansion techniques: Cossim with CCV, Cossim with TF-IDF, Cossim with Doc2Vec, SVM with Doc2Vec, Cossim with SBERT, SVM with SBERT, and SciDr-out.

5.1.4 Performance Metrics Our benchmark uses the following metrics:

Out-expansion accuracy: is the accuracy of predicting the right expansion for a given acronym in a textual document. Intuitively, this is the fraction of acronym-expansions that are correctly predicted. This corresponds to a micro-average of Precision and Recall, but, since we always predict an expansion for some acronym in the out-expansion task, those two metrics are equal. Accuracy is also used in previous out-expansion works [16, 42, 73] and analogous benchmarks,

e.g., for Word-Sense-Disambiguation [60]. Note that an acronym may appear many times in the same document and many times across documents. In our measure, if *A* is in *k* documents, it is counted *k* times, but if *A* is present *j* times in the same document, it is counted only once in that document.

Thus, for a test set of documents *D*, the out-expansion accuracy is: $\frac{\sum_{d \in D} |\text{correct distinct expansions for } d|}{\sum_{d \in D} |\text{distinct acronyms inside } d|}$.

Out-Expansion macro averages: Recently, Veyseh et al. [74][76] started using a different set of metrics that we have implemented and measured for completeness. Those metrics are macro-averages of Precision, Recall and F1-measures for acronym-expansions pairs. So, we calculate precision, recall, and F1-measure independently for each acronym-expansion in the training data. Then, averages of those measures are performed in order to obtain the final macro averages. Note that, when using these measures, very rare acronyms and expansions in the dataset will have the same impact as more frequent expansions and acronyms.

Representator execution time: is the execution time to create representations of training documents.

Average execution time per document: is the average execution time to predict expansions for acronyms in a document.

5.2 Out-expansion Experimental Results

Setup. For out-expansion on the benchmark presented in Section 5.1, we ran the experiments on a GoogleCloud platform¹³ machine with the following specifications: Intel Broadwell CPU platform with 8 cores, 30GB to 80GB of RAM (Random Access Memory). For SBERT, MadDog-out, SciDr-out, and LUKE half of a Tesla K80 GPU board was used. The code ran in Python 3.7.

To reduce the duration of experiments, we first find the representator's hyperparameters with the cosine similarity because it involves no learning nor hyperparameters of its own, then save the best representator model on disk. Next, given the best representator hyperparameters, we find the best out-expansion predictor model hyperparameters.

Results. For each dataset, in Table 8, we report the out-expansion accuracy and macro F1-measure to predict the expansions of acronyms in a document. In Table 9, we present the average execution times that out-expansion techniques that to process a document for each dataset. The *Technique Group* column identifies the out-expansion group that the technique belongs to, as organized in Section 5.1.3 (e.g., Classical). The *Predictors* column identifies the out-expansion predictor technique (e.g., Cossim or an ML classifier) that takes a given document representation to predict an expansion (e.g., Cossim). The *Representators* column indicates the technique used to generate a document representation (e.g., Doc2Vec). We did not run SciDr-out with External Data on CSWiki dataset because the external data (i.e., Wikipedia data) would overlap with CSWiki itself. The execution time of each ensemble technique is just the additional time required to decide on an expansion given the input predictions and confidence measures.

¹¹https://www.sbert.net/docs/pretrained_models.html#model-overview

¹²<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

¹³<https://cloud.google.com/>

Out-expansion Technique			ScienceWISE		MSH		CSWiki		SciAD		Average	
Technique Group	Predictors	Representators	Acc	MaF1	Acc	MaF1	Acc	MaF1	Acc	MaF1	Acc	MaF1
Classical	Random		47.72%	46.10%	47.04%	45.49%	14.54%	14.21%	33.06%	32.22%	35.59%	34.51%
	Most Frequent		70.52%	49.31%	50.30%	32.32%	47.76%	20.37%	69.08%	37.64%	59.41%	34.91%
	Cossim	CCV	91.34%	80.72%	97.62%	97.65%	77.96%	65.01%	92.28%	89.03%	89.80%	83.10%
		DCV	89.51%	78.69%	96.18%	96.07%	78.59%	65.86%	93.67%	87.08%	89.49%	81.92%
		SBE	88.07%	75.89%	95.84%	95.24%	74.60%	63.30%	86.50%	80.76%	86.25%	78.80%
Thakker		87.77%	77.38%	92.53%	91.68%	73.16%	63.86%	84.36%	73.21%	84.46%	76.53%	
Entity Disam.	LUKE		83.42%	57.33%	67.47%	58.95%	52.65%	46.60%	50.68%	42.53%	63.55%	51.35%
Sentence-Oriented	UAD		43.69%	46.73%	93.92%	92.55%	12.94%	11.60%	34.98%	45.75%	46.38%	49.16%
	MadDog-out		89.13%	68.84%	94.09%	93.16%	57.03%	47.71%	87.38%	73.23%	81.91%	70.73%
	SciDr-out		88.22%	77.45%	97.23%	96.76%	84.19%	72.67%	94.48%	88.94%	91.03%	83.96%
	SciDr-out with External Data		89.89%	77.86%	97.58%	97.22%	N/A	N/A	94.71%	89.42%	N/A	N/A
Representator	Cossim	TF-IDF	91.26%	81.82%	97.62%	97.57%	77.80%	65.36%	91.79%	83.48%	89.62%	82.06%
		LDA	85.56%	73.94%	93.81%	93.28%	71.89%	60.49%	84.56%	73.39%	83.95%	75.28%
		Doc2Vec	92.86%	83.14%	98.33%	98.07%	77.16%	65.25%	92.05%	82.96%	90.10%	82.35%
		SBERT	94.83%	85.32%	98.78%	98.80%	81.47%	67.67%	94.19%	89.76%	92.32%	85.39%
Classification	RF	TFIDF	70.82%	52.03%	84.53%	76.78%	32.11%	23.14%	87.64%	68.90%	68.77%	55.21%
		LDA	70.75%	54.13%	95.64%	92.84%	67.57%	50.78%	82.32%	61.33%	79.07%	64.77%
		Doc2Vec	79.18%	61.34%	96.58%	95.37%	66.29%	41.55%	84.39%	62.43%	81.61%	65.17%
	LR	TFIDF	71.05%	54.47%	93.41%	88.29%	71.89%	45.59%	80.63%	55.06%	79.24%	60.85%
		LDA	71.13%	51.49%	88.66%	80.02%	71.73%	48.77%	80.08%	55.16%	77.90%	58.86%
		Doc2Vec	88.83%	78.35%	98.87%	98.72%	76.68%	57.97%	90.75%	77.95%	88.78%	78.25%
	SVM	TFIDF	81.84%	62.13%	94.71%	91.27%	77.16%	53.54%	91.01%	78.24%	86.18%	71.29%
		LDA	78.88%	59.80%	93.64%	91.16%	71.89%	51.11%	85.59%	70.63%	82.50%	68.18%
		Doc2Vec	89.67%	79.31%	98.93%	98.79%	77.00%	58.70%	91.56%	80.88%	89.29%	79.42%
		SBERT	93.01%	83.91%	98.87%	98.84%	82.43%	64.44%	92.34%	86.53%	91.66%	83.43%
Combination of Representators	Cossim	CCV + Doc2Vec	90.27%	79.04%	98.19%	97.95%	77.16%	65.25%	86.92%	82.82%	88.14%	81.27%
		DCV + Doc2Vec	90.27%	79.01%	98.33%	98.10%	77.16%	65.19%	92.05%	82.96%	89.45%	81.32%
	SVM	CCV + Doc2Vec	89.97%	80.44%	98.95%	98.84%	77.00%	58.70%	80.73%	76.06%	86.66%	78.51%
		DCV + Doc2Vec	89.67%	79.35%	98.95%	98.83%	77.00%	58.70%	90.20%	75.90%	88.95%	78.20%
Ensemblers	Hard		94.15%	86.95%	99.60%	99.58%	84.19%	78.32%	96.59%	91.88%	93.63%	89.18%
	Soft		93.62%	85.00%	99.38%	99.41%	86.26%	79.33%	95.94%	91.04%	93.80%	88.70%

Table 8: Out-expansion accuracy (Acc) and macro F1-measure (MaF1). Values marked as bold indicate the best Acc obtained by an individual technique and by an ensembler, respectively in that dataset. A technique T1 is considered better than T2 if a non-parametric significance test (based on shuffling[27]) indicates that the difference in their means has a p-value < 0.05. Thus, even though each column has a highest mean value for some technique H which will be bolded, the value of a technique T will also be bolded if H is no better than T based on the p-value criterion. We apply the same p-value criteria to bold ensemblers on all datasets, except on ScienceWISE where we apply the statistical test to each ensembler against Cossim SBERT (the best technique on ScienceWISE).

In these out-expansion experiments, we measure the accuracy and macro F1 only on the acronym-expansions pairs whose acronym is ambiguous (i.e., have at least two expansions in the training data) and whose in-expansions are in the training data.

The best individual techniques (average above 89% of accuracy) in descending order are: Cossim with SBERT, SVM with SBERT, SciDr-out, Cossim with CCV, Cossim with TF-IDF, Cossim with DCV, Cossim with Doc2Vec alone or with DCV, and SVM with Doc2Vec. Regarding statistical significance, Cossim with SBERT

is the best for SciWISE. For MSH, SVM with Doc2vec combined with either CCV or DCV score higher accuracy. However, they are not statistically significantly better than: SVM with either Doc2Vec or SBERT, Cossim with SBERT, and LR with Doc2Vec. SciDr-out achieves higher accuracy for CSWiki, but is not statistically better than SVM with SBERT. Finally, for SciAD, SciDr-out with external data scores higher accuracy but not statistically significantly better than: SciDr-out and Cossim with SBERT.

Out-expansion Technique		Execution Times per Document (s)						
Technique Group	Predictors	Representators	SciWISE	MSH	CSWiki	SciAD	Average	
Classical	Random		0.00	0.00	0.00	0.00	0.00	
	Most Frequent		0.00	0.00	0.00	0.00	0.00	
	Cossim	CCV		0.01	0.06	0.07	0.12	0.07
		DCV		0.02	0.18	0.08	0.34	0.15
		SBE		0.02	0.07	0.06	0.06	0.05
Thakker		2.19	7.56	2.62	2.78	3.79		
Entity Disam.	LUKE		0.88	3.38	63.74	0.31	17.07	
Sentence-Oriented	UAD		0.00	0.01	0.01	0.00	0.01	
	MadDog-out		0.07	0.22	1.15	0.04	0.37	
	SciDr-out		0.72	1.19	2.51	0.70	1.28	
	SciDr-out with External Data		2.15	3.37	N/A	1.28	N/A	
Representator	Cossim	TF-IDF	0.06	5.02	4.20	0.87	2.53	
		LDA	0.01	0.02	0.03	0.02	0.02	
		Doc2Vec	0.01	0.01	0.36	0.02	0.10	
		SBERT	0.19	0.39	0.50	0.13	0.30	
Classification	RF	TFIDF	1.03	17.77	72.85	3.52	23.79	
		LDA	1.01	1.45	1.02	1.32	1.20	
		Doc2Vec	0.12	1.44	2.14	1.75	1.36	
	LR	TFIDF	0.06	5.03	41.36	0.92	11.84	
		LDA	0.01	0.02	0.03	0.02	0.02	
		Doc2Vec	0.01	0.02	0.36	0.04	0.11	
	SVM	TFIDF	0.04	4.77	7.05	1.05	3.23	
		LDA	0.01	0.03	0.03	0.02	0.02	
		Doc2Vec	0.01	0.01	0.37	0.02	0.10	
		SBERT	0.14	0.26	0.59	0.17	0.29	
Combination of Representators	Cossim	CCV + Doc2Vec	0.03	0.23	0.65	0.41	0.33	
		DCV + Doc2Vec	0.44	2.32	1.29	2.54	1.65	
	SVM	CCV + Doc2Vec	0.04	0.29	0.78	0.46	0.39	
		DCV + Doc2Vec	0.47	2.36	1.37	2.51	1.68	
Ensemblers	Hard		0.00	0.00	0.00	0.00	0.00	
	Soft		0.00	0.00	0.00	0.00	0.00	

Table 9: Overall out-expansion techniques average execution times per document.

Representators	SciWISE	MSH	CSWiki	SciAD	Average
CCV	0	1	5	1	2
DCV	0	1	5	1	2
SBE	6	23	115	9	38
LUKE	1 917	9 448	29 319	9 435	12 529
UAD	43	93	331	54	130
TF-IDF	2	10	58	1	18
LDA	155	7 766	8 830	4 247	5 250
Doc2Vec	13	32	212	108	91
SBERT	111	437	1 860	280	672
SciDr-in	11 227	42 716	147 454	50 282	62 920
SciDr-in External Data	14 299	46 651	N/A	66 441	42 464
MadDog-out	566	728	10 008	1 198	3 125

Table 10: Representator execution times in seconds for each dataset.

Interpretation: An important question in interpreting these numerical results is to understand why some techniques are better than others.

For out-expansion, the best approaches SciDr-out and Cossim/SVM with SBERT are based on language models trained on large data collections, but that does not tell the whole story. SciDr-out uses the particularly effective strategy of predicting the expansion span from the list of possible expansions passed as input. Further, SciDr-out is an ensemble of models trained in a 5-fold cross-validation setting. SBERT augments transformer language models to sentence similarity tasks using a siamese architecture that generates embeddings for each sentence and is trained to maximize similarity. Those embeddings turn out to be very informative regarding the context for documents: both Cossim or SVM combined with SBERT obtained on average the highest accuracy among individual techniques.

While LUKE’s transformer language model enables the creation of entity embeddings (in LUKE’s case, fine-tuned for expansion embeddings), the results are not the best for acronyms, even with fine-tuning. One reason is that each entity is referenced frequently (over 600 times on the average [81]). Acronym/expansion pairs are referenced less than twice on the average. For example, Wikipedia

Out-expansion Technique		Sci WISE	MSH	CS Wiki	Sci AD	Average
Predictors	Representators					
Cossim	CCV	0.35	0.25	0.38	0.34	0.33
	TF-IDF	0.37	0.33	0.48	0.30	0.37
	Doc2Vec	0.27	0.14	0.32	0.19	0.23
	SBERT	0.24	0.15	0.32	0.24	0.24
SVM	Doc2Vec	0.10	0.03	0.36	0.19	0.17
	SBERT	0.07	0.06	0.34	0.23	0.18

Table 11: Pearson correlation coefficient values of confidence with correct acronym expansion for each dataset and each technique. Accuracy is correlated with confidence, but only modestly.

has 11 million entity occurrences and 18 thousand distinct entities [81], an average of 611 mentions per entity. By contrast, examples of sentences whose acronyms are expanded by links are limited, because some acronyms may be defined in a single document. On CSWiki, for example, we have 10,400 acronym expansion occurrences (without counting repeated occurrences in the same document) and 8,600 distinct acronym/expansion pairs, an average of only 1.2 occurrences per acronym expansion.

Independently of which technique is best, we should note that each of the top techniques, except SciDr-out, gives a confidence score. For some of the best techniques SBERT, Doc2Vec, TFIDF, and CCV, the confidence score has a positive correlation with accuracy, though the correlation is modest (under 0.5), as indicated in Table 11. This low positive correlation is reflected in our results for ensemble techniques. The soft ensemble technique (in which each underlying technique’s weight is monotonic with its confidence) does well thanks to the positive correlation. On the other hand, hard voting ensemble techniques (in which each underlying technique votes for its preferred expansion regardless of confidence) perform even better, suggesting that the "wisdom of crowds" effect is stronger than using confidences. A deeper look at ensemble techniques for acronym expansion is a subject for future work.

Analysis of Classical techniques. The baselines Random and Most frequent out-expanders have better accuracies on SciWISE (48% and 71%) followed by SciAD (33% and 69%) and MSH (47% and 50%). The worst scores are obtained on CSWiki (15% and 48%), because that dataset has far more ambiguity.

MadDog-out achieves near the best performance on SciWISE and MSH, is also good on SciAD, and is worst on CSWiki. MadDog-out’s input is a sentence at a time like SBE and UAD, so context is also limited to a few words. In contrast, the best techniques, except SciDr-out, process all words in a document.

Analysis of remaining Representators and Classification techniques. LDA works well in many Natural Language Processing tasks, but less well for our out-expansion task, probably because measuring document similarity by comparing topics is not the best use of LDA.

Techniques with TF-IDF as a representator are surprisingly competitive in SciWISE and MSH, but they achieve inferior accuracy in CSWiki. Pre-processing steps (Section 5.1.2) play an important role in CCV performance, for instance, in the Tokenization step, we

remove stop-words which TF-IDF takes into account thus automatically giving them less relevance due to the IDF score.

Predictors using RF are slower and slightly less accurate than Cossim and SVMs. RFs struggle with the fact that, in this acronym expansion, there are many features but only a small number of samples, e.g., 300 dimensions from Doc2Vec and a few documents per acronym. By contrast, Cossim and SVM are usually the best predictors closely followed by LR. The representator hyperparameter search was performed for TF-IDF, LDA, and Doc2Vec using Cossim as a predictor, hence they are fitted for this predictor. SVMs and LRs are similar classifiers, the first fits a separating hyper-plane, the other a logistic function.

Execution times of representators. The training execution times depend only on the representators and are reported in Table 10. The Doc2Vec models used by Thakker et al. are created during the input document processing hence are not reported in Table 10.

The CCV and DCV representators take the least time (average 2s) closely followed by TF-IDF (average 18s) (Table 10). By contrast, Word2Vec (SBE and UAD) and Doc2Vec models take more time depending on the hyperparameters and dataset size (6-331s). LDA takes on average 5ks. The most expensive models are SciDr-out (14ks-66ks) followed by LUKE (1Ks-13ks) and MadDog-out (566s-10ks) which use either language models or neural networks. Thakker et al. [73] does not report execution time to produce the representator model, because the system builds a Doc2Vec per acronym during the input document processing time.

Document processing execution times. As observed in Table 9, among these best techniques, Cossim with CCV is the fastest for all datasets, able to process input documents in less than 0.07 seconds on dataset average. However, SVM with Doc2Vec is the fastest for MSH and SciWISE. The slowest among the best is Cossim with TF-IDF (average 2.5s), followed by SciDr-out (1.3s for base and 2.3s with external data). These differences are statistically significant.

External data analysis. SciDr-out with External Data improves over the SciDr-out base (Table 8), indicating that adding additional models trained on external data usually helps at roughly double the time cost. The downside is that training time more than doubles (we have to sum the SciDr-out External data in Table 10). Average execution times (Table 8) to process a document almost doubles as well (from 1.3s to 2.3s).

In summary:

- If neither training time nor document processing time is of major concern and especially if GPU processing is available, then use either a **Hard** ensembler (best but slowest), **SciDr-out** (best with more domain data) or **Cossim/SVM** with **SBERT** (fastest and close to best).
- A pipeline that balances time and accuracy is to use **Doc2vec** as feature inputs for either **Cossim** or **SVMs**.
- If training and test time is limited, use **Cossim** with **CCV**, which requires almost no training time (less than 5s) and is the fastest in testing time among the best set of techniques.

6 End-to-end Benchmark and Evaluation

The end-to-end benchmark described in Section 6.1 uses input documents containing acronyms, where some of those acronyms contain expansions and others do not. The evaluation in Section 6.2

will measure the recall and precision of acronym expansion for both student annotators and AcX pipelines.

6.1 A Benchmark of End-to-End Acronym Expansion

Section 6.1.1 describes the train and test datasets used in this benchmark. Section 6.1.2 lists the end-to-end acronym expander systems included in the benchmark. Finally, Section 6.1.3 presents the metrics that our benchmark uses to evaluate the systems.

6.1.1 Datasets The end-to-end benchmark uses two different datasets: (i) for testing, the end-to-end dataset of Section 4.1.1. (ii) The *train* dataset consists of documents from Wikipedia that do not belong to the annotated test set. Those documents came from the Wikipedia dump of March 1, 2020¹⁴. These were converted to pure text using WikiExtractor [4].

We preprocessed all the documents using all the steps described in Section 5.1.2 for all out-expansion techniques except MadDog-out which uses its own preprocessing techniques.

6.1.2 End-to-end systems We use: (i) the end-to-end MadDog System (**MadDog-sys**) and (ii) various **pipelines of AcX** consisting of an in-expansion technique along with possibly the Link Follower (**LF**) technique followed by an out-expansion technique possibly with machine learning (see Figure 1). An example of a pipeline would be the SH in-expander, followed by the LF component, Doc2Vec, and SVMs. The pipelines we test consist of combinations of the most practical (accurate and fastest) techniques for in-expansion and out-expansion as determined by the benchmarks in Sections 4.2 and 5.2. Specifically, AcX pipelines use either the **MadDog-in** or the **SH** technique as in-expanders to identify acronyms and expansions in input documents. For out-expansion, AcX pipelines include one of the following combinations of out-expansion techniques, i.e., a predictor (Section 3.3) with a representator (Section 3.2): (i) **Cossim** with **SBERT**; (ii) **SVM** with **SBERT**; (iii) **Cossim** with **CCV**; (iv) **Cossim** with **Doc2vec**; and (v) **SVM** with **Doc2vec**.

Finally, we compare the accuracies of MadDog-sys, the various pipelines of AcX, and the **student annotators** (before the reviewer, the third annotator, resolved conflicts to decide on the final annotations).

6.1.3 Performance Metrics Similarly to Section 4.1.3, we evaluate MadDog-sys, different pipelines of AcX, and human annotators listed in Section 6.1.2 in terms of Precision (**P**), Recall (**R**) and F1-Measure (**F1**). *Precision* is the number of correct system-found acronym-expansion pairs divided by the total number of system-identified pairs. *Recall* is the number of correct system-found acronym-expansion pairs divided by the total number of human-found pairs. *F1-measure* is the harmonic mean of Precision and Recall. We also measure training and per test document execution times.

6.2 Results on End-to-end Experiments

Setup. For these experiments, we used a virtual machine with the following specifications: AMD EPYC Processor with 16 cores and 256GB of RAM (Random Access Memory). For SBERT, the virtual machine specifications were: five cores of an Intel Xeon Gold 6126 Processor, 40GB of RAM and a NVIDIA GeForce RTX 2080 Ti. The code ran in Python 3.7.

Results. Table 12 presents the results for the AcX system running each one of the different pipelines mentioned in Section 6.1.2, the MadDog-sys¹⁵, and the results for the student annotators. Moreover, apart from the end-to-end quality metrics in Table 12, we also report the quality metrics for out-expansion (i.e., acronyms left to expand after in-expansion and LF component).

The AcX pipeline composed by MadDog-in, SVM with SBERT without Link Following (LF) obtains the best results with precision (61.32%) and F1-measure (54.97%). However, based on the F1-measure, this is not statistically significantly better (i.e., P-value above 0.05) than the following system pipelines: (i) with the same out-expander but with LF, or (ii) with SH and SVM with SBERT, without LF. The best system pipeline takes **2s** on average to process a document. Our best AcX pipeline obtains better results for all measures than the MadDog-sys (**+20%** of F1) and is faster (**2s** to **1084s**).

Best AcX pipeline analysis. The reason why the precision is low is that our best AcX pipeline considers certain strings to be acronyms even though they are not (245 in total). Some are small words like "and" and "not". Others are codes like ZAB and ZAU that refer to airports. Conversely, the acronym and in-expansion extraction component fails to identify lower case acronyms as acronyms, common measurement units (e.g., m for meter, g for gram, kbit for kilobit) and some common language abbreviations (e.g., Micro, "etc", email) which usually everyone knows. By contrast, AcX provides the correct expansion for the acronyms that newcomers to a field may not know, e.g., CAS - Computer Algebra System; SLS - SoftLanding Linux system; and ILM - Industrial Light & Magic.

In and out expansion analysis. We report independently the performance of in-expansion in Table 13 and out-expansion in Table 12. When evaluating just the acronym and in-expander extraction component of AcX pipelines, using SH scored an in-expansion F1-measure of **68.88%** and using MadDog-in scored **69.92%**. If we evaluate out-expansion (acronyms left to expand after in-expansion and LF component), our best AcX pipeline (SVM as predictor with SBERT as representator) obtains an F1-measure of **47.76%**.

LF analysis. Overall, the Link Following component improves the quality metrics marginally if at all for these link-poor datasets. Unsurprisingly, the better the out-expander, the less LF helps.

We report Precision for only the LF predictions in Table 13: **77.66%** when running with SH and **73.19%** when running with MadDog-in. When combined with the best out-expander techniques (i.e., SVM as predictor with SBERT as representator), following links reduces the score by **-0.52%**, however when we consider the second-best set of out-expander techniques (i.e., SVM as predictor with Doc2Vec as representator), LF improves the F1-score by **+0.37%**

¹⁴<https://dumps.wikimedia.org/enwiki>

¹⁵<https://archive.org/details/MadDog-models>

AcX (pipelines)				Out-expansion			End-to-end			Execution Time per Document (s)		
In-expander	Out-expander Predictor	Representator	LF	P	R	F1	P	R	F1			
SH	Cossim	CCV	No	44.00%	37.40%	40.44%	51.43%	45.62%	48.35%	21.31		
			Yes	34.66%	22.68%	27.42%	52.16%	46.30%	49.05%	18.51		
		SBERT	No	46.03%	39.12%	42.30%	53.10%	47.10%	49.92%	0.15		
	SVM	Doc2Vec	Yes	36.70%	23.99%	29.01%	53.34%	47.35%	50.16%	1.16		
			Yes	40.86%	26.74%	32.33%	56.05%	49.75%	52.71%	3.26		
		SBERT	No	51.10%	43.43%	46.95%	57.27%	50.80%	53.84%	2.37		
MadDog-in	Cossim	CCV	Yes	42.02%	27.46%	33.22%	56.68%	50.31%	53.30%	2.91		
			No	44.77%	34.62%	39.05%	53.12%	43.15%	47.62%	17.45		
		Doc2Vec	Yes	37.36%	21.85%	27.57%	54.46%	44.44%	48.95%	19.51		
			Yes	40.20%	23.50%	29.66%	56.20%	45.86%	50.51%	7.73		
		SBERT	No	57.27%	36.55%	41.23%	55.17%	44.81%	49.46%	0.33		
			Yes	40.02%	23.29%	29.44%	56.28%	45.93%	50.58%	7.76		
	SVM	Doc2Vec	No	52.08%	40.27%	45.42%	59.12%	48.02%	53.00%	1.12		
			Yes	45.38%	26.53%	33.49%	59.38%	48.46%	53.37%	7.63		
		SBERT	No	54.76%	42.35%	47.76%	61.32%	49.81%	54.97%	2.29		
			Yes	47.09%	27.40%	34.64%	60.59%	49.44%	54.45%	9.12		
		MadDog-sys				25.40%	18.38%	21.33%	37.85%	29.14%	32.93%	1084.92
		Student annotators				N/A	N/A	N/A	88.36%	76.41%	81.95%	N/A

Table 12: Out-expansion and end-to-end system quality metrics and average execution times to process a document in seconds. The in-expander technique and the Link Follower (LF) component provide the input for out-expansion techniques, so when those two components are not fixed, we cannot compare out-expansion techniques directly using the out-expansion quality metrics. End-to-end values marked as bold indicate the best obtained in that metric. A method M1 is considered better than M2 if a non-parametric significance test (based on shuffling[27]) indicates that the difference in their means has a p-value < 0.05. Thus, even though each column has a highest mean value for some method H, the value of a method M will be bolded if H is no better than M based on the p-value criterion.

In-expander Technique	In-expansion			LF
	P	R	F1	P
SH	86.17%	57.37%	68.88%	77.66%
MadDog-in	91.49%	56.58%	69.92%	73.19%

Table 13: In-expansion prediction, recall, and F1-measure and Link Follower (LF) prediction when each in-expander technique is used.

Out-Expansion Representator	In-expansion Technique		
	SH	MadDog-in	Average
CCV	389	332	361
Doc2Vec	15,985	13,341	14,663
SBERT	42,769	33,724	38,247

Table 14: Representator execution times in seconds for each in-expander technique used to process the train dataset.

(Table 12). Thus, the better the out-expander, the less link following helps.

For the link follower, execution times were not measured for the time to download a linked document. When we predict an acronym with LF then it is an acronym less to out-expand. For SVM with SBERT or Doc2Vec we have increments of 6-7s by following links with MadDog-in, however for Cossim with CCV with MadDog-in0 and follow links with SH we observe a small difference. If the out-expansion technique is slow, then following links may also improve execution times.

Execution times of representators. We report in Table 14, the execution times to create the out-expansion representators for each system pipeline. We observe that CCV executes in less time (389s with SH and 332s with MadDog-in) than Doc2Vec (15,985s with SH and 13,341s with MadDog-in). SBERT representator takes the longest time (42,769 with SH and 33,724s with MadDog-in)

Comparison with human performance. Compared with human annotators, our best AcX pipeline (MadDog-in and SVMs with SBERT) is around 27% lower in Precision, Recall, and F1-Measure. So, there is a lot of room for improvement. On the other hand, automatic Acronym Expansion is rapid (2s per document) and can give at least a good first guess.

An example application of AcX. Consider one of the documents out of the 163 at random whose original page is here [https://en.wikipedia.org/wiki/CC_\(complexity\)](https://en.wikipedia.org/wiki/CC_(complexity)). Our best AcX pipeline identified the following acronym-expansion pairs: CC - comparator circuits; CCVP - comparator circuit value problem; AC - alternating current; NC - nick's class; and NL - national league. However, it failed to identify CC-complete, and P. We can see that CC, CCVP, NC, and AL are correct and NL is incorrect. With a different Pipeline consisting of Doc2Vec instead of SBERT, AL is incorrect, but NL is correct.

In summary:

- The best AcX pipeline consists of **MadDog-in**, with **SVM** and **SBERT**.
- The **Link Following** component has high precision but does not improve AcX performance compared with the best pipeline, though that could change depending on the density of acronym-related links.

7 Conclusions and Future Work

The AcX system synthesizes and extends the best of previous work on acronym expansion. In the process, our major technical findings are:

- In-expansion rule-based techniques (SH and MadDog-in) usually work best and require little execution time.
- For out-expansion, SciDr-out and Cossim or SVMs with SBERT usually work best, followed by Cossim and SVMs with either CCV or Doc2Vec.
- There is still a significant gap between the best AcX pipelines and human-level performance.

There are five data and software products of our work that future researchers can either extend or use as a basis of comparison.

- (1) The first human-annotated dataset for end-to-end acronym expander systems.
- (2) Three benchmarks to evaluate: (i) in-expansion techniques, (ii) out-expansion techniques, (iii) the combination in an end-to-end setting.
- (3) The end-to-end AcX system is available publicly and can be applied to arbitrary languages, follows hyperlinks, and can incorporate new in- and out-expansion techniques.

Future Work

Because the automated techniques in the state-of-the-art fall well below human-level accuracy levels, there is a large margin for improvement. We see the need for improvements in both in-expansion (especially acronym identification) and out-expansion. Some promising avenues for improvements include: (i) more accurate in-expansion (e.g., additional acronym-expansion extraction patterns), (ii) new context representation techniques, and (iii) an extensive study of ensemble techniques.

With respect to the AcX system, we will add an Application Programming Interface (API) so text analytics systems (e.g., entity disambiguation or sentiment analysis) can benefit from acronym expansion. Finally, because our platform easily extends to other languages (e.g., our Portuguese extension was done by a high school student), we plan to create AcX pipelines for a variety of natural languages.

Acknowledgments

Pereira's work was supported by national funds through FCT (*Fundação para a Ciência e a Tecnologia*), under the PhD Scholarship SFRH/BD/135719/2018. Furthermore, Pereira and Galhardas' work was supported by national funds through FCT under the project UIDB/50021/2020.

Shasha's work has been partly supported by (i) the New York University Abu Dhabi Center for Interacting Urban Networks (CITIES), funded by Tamkeen under the NYUAD Research Institute Award CG001 and by the Swiss Re Institute under the Quantum Cities initiative, (ii) NYU WIRELESS, (iii) U.S. National Science Foundation grants 1934388, 1840761, and 1339362, and (iv) INRIA.

The server virtual machine used to run the experiments was supported by BioData.pt – *Infraestrutura Portuguesa de Dados Biológicos*, project 22231/01/SAICT/2016, funded by Portugal 2020. This material is based upon work supported by Google Cloud. This

work was carried out on the Dutch national e-infrastructure with the support of SURF Cooperative.

Finally, we would like to thank the reviewers for several excellent suggestions.

References

- [1] ABBREX. 2011. ABBREX - The Abbreviation Expander. <http://abbrex.com/>
- [2] Khaled Abdalgader and Andrew Skabar. 2012. Unsupervised Similarity-based Word Sense Disambiguation Using Context Vectors and Sentential Word Importance. *ACM Transactions on Speech and Language Processing* 9, 1 (2012), 2–21.
- [3] Hiroko Ao and Toshihisa Takagi. 2005. ALICE: an algorithm to extract abbreviations from MEDLINE. *Journal of the American Medical Informatics Association* 12, 5 (2005), 576–586.
- [4] Giuseppe Attardi. 2015. WikiExtractor. <https://github.com/attardi/wikiextractor>.
- [5] S Azimi, H Veisi, and R Amouie. 2019. A method for automatic detection of acronyms in texts and building a dataset for acronym disambiguation. In *Iranian Conference on Signal Processing and Intelligent Systems*. 1–4.
- [6] Adrian Barnett and Zoe Doubleday. 2020. Meta-Research: The growth of acronyms in the scientific literature. *eLife* 9 (jul 2020), e60080. <https://doi.org/10.7554/eLife.60080>
- [7] Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: Pretrained Language Model for Scientific Text. In *Empirical Methods in Natural Language Processing*.
- [8] Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural Language Processing with Python* (1st ed.). O'Reilly Media, Inc.
- [9] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent Dirichlet Allocation. *Journal of Machine Learning Research* 3 (March 2003), 993–1022.
- [10] Robert Bossy, Louise Deléger, Estelle Chaix, Mouhamadou Ba, and Claire Nédellec. 2019. Bacteria Biotope at BioNLP Open Shared Tasks 2019. In *Workshop on Biomedical Natural Language Processing Open Shared Tasks*. Association for Computational Linguistics, Hong Kong, China, 121–131. <https://doi.org/10.18653/v1/D19-5719>
- [11] Leo Breiman. 2001. Random Forests. *Machine Learning* 45, 1 (2001), 5–32.
- [12] Jean Charbonnier and Christian Wartena. 2018. Using Word Embeddings for Unsupervised Acronym Disambiguation. In *International Conference on Computational Linguistics*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2610–2619. <https://www.aclweb.org/anthology/C18-1221>
- [13] Zheng Chen, Suzanne R Tamang, Adam Lee, Xiang Li, Wen-Pin Lin, Matthew G Snover, Javier Artilles, Marissa Passantino, and Heng Ji. 2010. CUNY-BLENDER TAC-KBP2010 Entity Linking and Slot Filling System Description. *Theory and Applications of Categories* (2010).
- [14] Daphné Chopard and Irena Spasić. 2019. A Deep Learning Approach to Self-expansion of Abbreviations Based on Morphology and Context Distance. In *Statistical Language and Speech Processing*. 71–82. https://doi.org/10.1007/978-3-030-31372-2_6
- [15] Manuel R. Ciosici and Ira Assent. 2018. Abbreviation Expander - a Web-based System for Easy Reading of Technical Documents. In *Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 1–4. <https://www.aclweb.org/anthology/C18-2001>
- [16] Manuel R. Ciosici, Tobias Sommer, and Ira Assent. 2019. Unsupervised Abbreviation Disambiguation Contextual disambiguation using word embeddings. *Computing Research Repository* arXiv:1904.00929 (2019). [arXiv:1904.00929](http://arxiv.org/abs/1904.00929) version 2.
- [17] Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20, 1 (1960), 37–46.
- [18] Nigel Collier and Jin-Dong Kim. 2004. Introduction to the Bio-entity Recognition Task at JNLPBA. In *International Joint Workshop on Natural Language Processing in Biomedicine and its Applications*. Geneva, Switzerland, 73–78. <https://aclanthology.org/W04-1213>
- [19] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20, 3 (1995), 273–297.
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Computing Research Repository* arXiv:1810.04805 (2018). <https://arxiv.org/abs/1810.04805>
- [21] Nicholas Egan and John Bohannon. 2021. Primer AI's Systems for Acronym Identification and Disambiguation. In *Workshop on Scientific Document Understanding co-located with 35th AAAI Conference on Artificial Intelligence*. CEUR-WS.org. <http://ceur-ws.org/Vol-2831/paper30.pdf>
- [22] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. 2008. LIBLINEAR: A Library for Large Linear Classification. *Journal of Machine Learning Research* 9 (jun 2008), 1871–1874.
- [23] Shicong Feng, Yuhong Xiong, Conglei Yao, Liwei Zheng, and Wei Liu. 2009. Acronym Extraction and Disambiguation in Large-Scale Organizational Web

- Pages. In *Conference on Information and Knowledge Management* (Hong Kong, China). Association for Computing Machinery, New York, NY, USA, 1693–1696. <https://doi.org/10.1145/1645953.1646206>
- [24] Michael R. Glass, Md. Faisal Mahbub Chowdhury, and Alfio Massimiliano Gliozzo. 2017. Language Independent Acquisition of Abbreviations. *Computing Research Repository* arXiv:1709.08074 (2017). arXiv:1709.08074 <http://arxiv.org/abs/1709.08074> version 1.
- [25] Aitor Gonzalez-Agirre, Montserrat Marimon, Ander Intxaurre, Obdulia Rabal, Marta Villegas, and Martin Kralligner. 2019. PharmaCoNER: Pharmacological Substances, Compounds and proteins Named Entity Recognition track. In *Workshop on Biomedical Natural Language Processing Open Shared Tasks*. Association for Computational Linguistics, Hong Kong, China, 1–10. <https://doi.org/10.18653/v1/d19-5701>
- [26] Phil Gooch. 2012. BADREX: In situ expansion and coreference of biomedical abbreviations using dynamic regular expressions. *Computing Research Repository* arXiv:1206.4522 (2012). arXiv:1206.4522 <http://arxiv.org/abs/1206.4522> version 1.
- [27] Phillip I Good. 2006. *Resampling Methods: A Practical Guide to Data Analysis*. Birkhäuser Basel. <https://doi.org/10.1007/0-8176-4444-X>
- [28] Richard D Hipp. 2020. SQLite. <https://www.sqlite.org/>
- [29] Rezarta Islamaj Doğan, Donald C Comeau, Lana Yeganova, and W John Wilbur. 2014. Finding abbreviations in biomedical literature: three BioC-compatible modules and four BioC-formatted corpora. *Database: the journal of biological databases and curation* 2014 (2014).
- [30] Paul Jaccard. 1912. The distribution of the flora in the alpine zone. 1. *New phytologist* 11, 2 (1912), 37–50.
- [31] Kayla Jacobs, Alon Itai, and Shuly Wintner. 2020. Acronyms: identification, expansion and disambiguation. *Annals of Mathematics and Artificial Intelligence* 88, 5 (2020), 517–532.
- [32] A Jain, S Cucercan, and Saliha Azzam. 2007. Acronym-Expansion Recognition and Ranking on the Web. *IEEE International Conference on Information Reuse and Integration* (2007), 209–214.
- [33] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2014. *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated.
- [34] Antonio J Jimeno-Yepes, Bridget T McInnes, and Alan R Aronson. 2011. Exploiting MeSH indexing in MEDLINE to generate a data set for word sense disambiguation. *BMC Bioinformatics* 12, 1 (2011), 1–14.
- [35] Karen Spärck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation* 28 (1972), 11–21.
- [36] İlknur Karadeniz, Ömer Faruk Tuna, and Arzuhan Özgür. 2019. BOUN-ISIK Participation: An Unsupervised Approach for the Named Entity Normalization and Relation Extraction of Bacteria Biotopes. In *Workshop on Biomedical Natural Language Processing Open Shared Tasks*. Association for Computational Linguistics, Hong Kong, China, 150–157. <https://doi.org/10.18653/v1/D19-5722>
- [37] Klaus Krippendorff. 2018. *Content analysis: An introduction to its methodology*. Sage publications.
- [38] Cheng-Ju Kuo, Maurice HT Ling, Woody Lin, and Chun-Nan Hsu. 2009. BIOADI: A machine learning approach to identifying abbreviations and definitions in biological literature. *BMC bioinformatics* 10 Suppl 15 (12 2009), S7.
- [39] Matt J Kusner, Yu Sun, Nicholas I Kolkin, and Kilian Q Weinberger. 2015. From Word Embeddings to Document Distances. In *International Conference on International Conference on Machine Learning*. JMLR.org, 957–966.
- [40] John D Lafferty, Andrew McCallum, and Fernando C N Pereira. 2001. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In *International Conference on Machine Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 282–289. <http://dl.acm.org/citation.cfm?id=645530.655813>
- [41] Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. In *International Conference on Machine Learning* (Beijing, China), Vol. 32. 1188–1196.
- [42] Chao Li, Lei Ji, and Jun Yan. 2015. Acronym Disambiguation Using Word Embedding. In *AAAI Conference on Artificial Intelligence* (Austin, Texas), 4178–4179.
- [43] Yang Li, Bo Zhao, Ariel Fuxman, and Fangbo Tao. 2018. Guess Me if You Can: Acronym Disambiguation for Enterprises. In *Annual Meeting of the Association for Computational Linguistics*, Vol. 1: Long Papers. Association for Computational Linguistics, Melbourne, Australia, 1308–1317. <https://doi.org/10.18653/v1/P18-1121>
- [44] Nicholas B Link, Sicong Huang, Tianrun Cai, Jiehuan Sun, Kumar Dahal, Lauren Costa, Kelly Cho, Katherine Liao, Tianxi Cai, and Chuan Hong. 2022. Binary acronym disambiguation in clinical notes from electronic health records with an application in computational phenotyping. *International Journal of Medical Informatics* 162 (2022), 104753. <https://doi.org/10.1016/j.ijmedinf.2022.104753>
- [45] Jie Liu, Caihua Liu, and Yalou Huang. 2017. Multi-granularity sequence labeling model for acronym expansion identification. *Information Sciences* 378 (2017), 462–474.
- [46] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. *Computing Research Repository* arXiv:1301.3781 (2013). <https://arxiv.org/abs/1301.3781> version 3.
- [47] Sungrim Moon, Bridget McInnes, and Genevieve B Melton. 2015. Challenges and practical approaches with word sense disambiguation of acronyms and abbreviations in the clinical domain. *Healthcare Informatics Research* 21, 1 (jan 2015), 35–42. <https://doi.org/10.4258/hir.2015.21.1.35>
- [48] Sungrim Moon, Serguei Pakhomov, and Genevieve B Melton. 2012. Automated disambiguation of acronyms and abbreviations in clinical texts: window and training size considerations. *AMIA Annual Symposium proceedings* 2012 (2012), 1310–1319.
- [49] Andrea Moro and Roberto Navigli. 2015. SemEval-2015 Task 13: Multilingual All-Words Sense Disambiguation and Entity Linking. In *International Workshop on Semantic Evaluation*. Association for Computational Linguistics, Denver, Colorado, 288–297. <https://doi.org/10.18653/v1/S15-2049>
- [50] Andrea Moro, Alessandro Raganato, and Roberto Navigli. 2014. Entity Linking meets Word Sense Disambiguation: a Unified Approach. *Transactions of the Association for Computational Linguistics* 2 (2014), 231–244. https://doi.org/10.1162/tacl_a_00179
- [51] Roberto Navigli. 2009. Word Sense Disambiguation: A Survey. *Comput. Surveys* 41, 2, Article 10 (Feb. 2009), 69 pages. <https://doi.org/10.1145/1459352.1459355>
- [52] Vincent Ng. 2017. Machine Learning for Entity Coreference Resolution: A Retrospective Look at Two Decades of Research. In *AAAI Conference on Artificial Intelligence*, Vol. 31. <https://ojs.aaai.org/index.php/AAAI/article/view/11149>
- [53] Serguei Pakhomov, Ted Pedersen, and Christopher G Chute. 2005. Abbreviation and acronym disambiguation in clinical discourse. *AMIA Annual Symposium proceedings* 2005 (2005), 589–593.
- [54] Youngja Park and Roy J. Byrd. 2001. Hybrid Text Mining for Finding Abbreviations and their Definitions. In *Empirical Methods in Natural Language Processing*. <https://www.aclweb.org/anthology/W01-0516>
- [55] Rebecca Passonneau. 2006. Measuring Agreement on Set-valued Items (MASI) for Semantic and Pragmatic Annotation. In *International Conference on Language Resources and Evaluation*. European Language Resources Association (ELRA), Genoa, Italy.
- [56] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [57] Anja Pilz and Gerhard Paaß. 2011. From Names to Entities Using Thematic Context Distance. In *Conference on Information and Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 857–866. <https://doi.org/10.1145/2063576.2063700>
- [58] Roman Prokofyev, Gianluca Demartini, Alexey Boyarsky, Oleg Ruchayskiy, and Philippe Cudré-Mauroux. 2013. Ontology-Based Word Sense Disambiguation for Scientific Literature. In *European Conference on Advances in Information Retrieval* (Moscow, Russia). Springer-Verlag, Berlin, Heidelberg, 594–605. https://doi.org/10.1007/978-3-642-36973-5_50
- [59] J Pustejovsky, J Castañó, B Cochran, M Kotecki, and M Morrell. 2001. Automatic extraction of acronym-meaning pairs from MEDLINE databases. *Studies in Health Technology and Informatics* 84, Pt 1 (2001), 371–375.
- [60] Alessandro Raganato, Jose Camacho-Collados, and Roberto Navigli. 2017. Word Sense Disambiguation: A Unified Evaluation Framework and Empirical Comparison. In *Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, Valencia, Spain, 99–110. <https://www.aclweb.org/anthology/E17-1010>
- [61] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Hong Kong, China, 3982–3992. <https://doi.org/10.18653/v1/D19-1410>
- [62] Saneesh Mohammed N and K A Abdul Nazeer. 2013. An improved method for extracting acronym-definition pairs from biomedical Literature. In *2013 International Conference on Control Communication and Computing (ICCC)*. 194–197.
- [63] Ariel S Schwartz and Marti A Hearst. 2003. A simple algorithm for identifying abbreviation definitions in biomedical text. In *Pacific Symposium on Biocomputing*. 451–462.
- [64] Özge Sevgili, Alexander Panchenko, and Chris Biemann. 2019. Improving Neural Entity Disambiguation with Graph Embeddings. In *Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*. Association for Computational Linguistics, Florence, Italy, 315–322. <https://doi.org/10.18653/v1/P19-2044>
- [65] Wei Shen, Yuhua Li, Yinan Liu, Jiawei Han, Jianyong Wang, and Xiaojie Yuan. 2021. Entity Linking Meets Deep Learning: Techniques and Solutions. *IEEE Transactions on Knowledge and Data Engineering* (2021). <https://doi.org/10.1109/TKDE.2021.3117715>
- [66] Wei Shen, Jianyong Wang, and Jiawei Han. 2015. Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions. *IEEE Transactions on Knowledge and Data Engineering* 27, 2 (2015), 443–460. <https://doi.org/10.1109/TKDE.2014.2327028>

- [67] Utpal Kumar Sikdar and Björn Gambäck. 2017. A Feature-based Ensemble Approach to Recognition of Emerging and Rare Named Entities. In *Workshop on Noisy User-generated Text*. Association for Computational Linguistics, Copenhagen, Denmark, 177–181. <https://doi.org/10.18653/v1/W17-4424>
- [68] Aadarsh Singh and Priyanshu Kumar. 2021. SciDr at SDU-2020: IDEAS–Identifying and Disambiguating Everyday Acronyms for Scientific Domain. In *Workshop on Scientific Document Understanding co-located with 35th AAAI Conference on Artificial Intelligence*. CEUR-WS.org. <http://ceur-ws.org/Vol-2831/paper31.pdf>
- [69] Sunghwan Sohn, Donald C Comeau, Won Kim, and W John Wilbur. 2008. Abbreviation definition identification based on automatic precision estimates. *BMC bioinformatics* 9, 1 (2008), 402.
- [70] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. MPNet: Masked and Permuted Pre-training for Language Understanding. *Advances in Neural Information Processing Systems* 33 (2020), 16857–16867.
- [71] Mark Stevenson, Yikun Guo, Abdulaziz Al Amri, and Robert Gaizauskas. 2009. Disambiguation of Biomedical Abbreviations. In *Workshop on Current Trends in Biomedical Natural Language Processing*. Association for Computational Linguistics, USA, 71–79.
- [72] Bilyana Taneva, Tao Cheng, Kaushik Chakrabarti, and Yeye He. 2013. Mining Acronym Expansions and Their Meanings Using Query Click Log. In *International Conference on World Wide Web*. Association for Computing Machinery, New York, NY, USA, 1261–1272. <https://doi.org/10.1145/2488388.2488498>
- [73] Aditya Thakker, Suhail Barot, and Sudhir Bagul. 2017. Acronym Disambiguation: A Domain Independent Approach. *Computing Research Repository* arXiv:1711.09271 (2017). <https://arxiv.org/abs/1711.09271> version 3.
- [74] Amir Pouran Ben Veyseh, Franck Deroncourt, Walter Chang, and Thien Huu Nguyen. 2021. MadDog: A Web-based System for Acronym Identification and Disambiguation. In *European Chapter of the Association for Computational Linguistics*.
- [75] Amir Pouran Ben Veyseh, Franck Deroncourt, Thien Huu Nguyen, Walter Chang, and Leo Anthony Celi. 2021. Acronym Identification and Disambiguation Shared Tasks for Scientific Document Understanding. In *Workshop on Scientific Document Understanding co-located with 35th AAAI Conference on Artificial Intelligence*. CEUR-WS.org. <http://ceur-ws.org/Vol-2831/paper33.pdf>
- [76] Amir Pouran Ben Veyseh, Franck Deroncourt, Quan Hung Tran, and Thien Huu Nguyen. 2020. What Does This Acronym Mean? Introducing a New Dataset for Acronym Identification and Disambiguation. In *International Conference on Computational Linguistics*.
- [77] Yonghui Wu, Joshua C Denny, S Trent Rosenbloom, Randolph A Miller, Dario A Giuse, Lulu Wang, Carmelo Blanquicett, Ergin Soysal, Jun Xu, and Hua Xu. 2017. 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* 24 (2017), 79–86. <https://doi.org/10.1093/jamia/ocw109>
- [78] Yonghui Wu, Jun Xu, Yaoyun Zhang, and Hua Xu. 2015. Clinical Abbreviation Disambiguation Using Neural Word Embeddings. In *Workshop on Biomedical Natural Language Processing*. Association for Computational Linguistics, Beijing, China, 171–176. <https://doi.org/10.18653/v1/W15-3822>
- [79] Vikas Yadav and Steven Bethard. 2018. A Survey on Recent Advances in Named Entity Recognition from Deep Learning models. In *International Conference on Computational Linguistics*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 2145–2158. <https://aclanthology.org/C18-1182>
- [80] Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. LUKE: Deep Contextualized Entity Representations with Entity-aware Self-attention. In *Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online, 6442–6454. <https://doi.org/10.18653/v1/2020.emnlp-main.523>
- [81] Ikuya Yamada, Koki Washio, Hiroyuki Shindo, and Yuji Matsumoto. 2019. Global Entity Disambiguation with BERT. *Computing Research Repository* arXiv:1909.00426 (2019). <https://arxiv.org/abs/1909.00426> version 3.
- [82] Anna Yarygina and Natalia Vassilieva. 2012. High-recall Extraction of Acronym-definition Pairs with Relevance Feedback. In *Joint Extending Database Technology and International Conference on Database Theory Workshops*. ACM, New York, NY, USA, 21–28. <https://doi.org/10.1145/2320765.2320781>
- [83] Hong Yu, Won Kim, Vasileios Hatzivassiloglou, and John Wilbur. 2006. A Large Scale, Corpus-Based Approach for Automatically Disambiguating Biomedical Abbreviations. *ACM Transactions on Information Systems* 24, 3 (jul 2006), 380–404. <https://doi.org/10.1145/1165774.1165778>
- [84] Danqing Zhu, Wangli Lin, Yang Zhang, Qiwei Zhong, Guanxiang Zeng, Weilin Wu, and Jiayu Tang. 2021. AT-BERT: Adversarial Training BERT for Acronym Identification Winning Solution for SDU@ AAAI-21. In *Workshop on Scientific Document Understanding co-located with 35th AAAI Conference on Artificial Intelligence*. CEUR-WS.org. <http://ceur-ws.org/Vol-2831/paper28.pdf>
- [85] Stefan Zwicklbauer, Christin Seifert, and Michael Granitzer. 2016. Robust and collective entity disambiguation through semantic embeddings. In *Special Interest Group in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 425–434. <https://doi.org/10.1145/2911451.2911535>

A Dataset creation and structures details

A.1 End-to-end dataset creation process

To further the annotation process, we created two forms: one that promoted the annotation task and allowed for volunteers to register¹⁶ and another one for the annotation process itself¹⁷. The first form explained the overall annotation task and stated the prize each participant could win. After filling in the first form, each participant automatically received an e-mail message with the Wikipedia documents to annotate and a link to the second form. Furthermore, in the second form additional details, instructions, and examples of the annotation process were given to each participant explaining the annotation process and how to fill in the form correctly.

A.2 Data structures used in AcX

Each dataset referenced in Section 4.1.1 and Section 5.1.1 were in different formats and used a different annotation notation for the acronym-expansion pairs in the dataset documents. To facilitate access to each dataset by the benchmark we created three main data structures from the original dataset:

- A dictionary that maps each document id to the corresponding raw text and for out-expansion an additional dictionary is provided with the preprocessed text.
- A dictionary that maps each document id to the acronym-expansion pairs whose acronyms are present in text.
- A dictionary that maps each acronym in the corpus to the corresponding in- and/or out- expansions with the document ids where they appear.

B In-expansion Glossary-level results

We report Glossary-level Precision, Recall, and F1-measure values for the biomedical datasets (i.e., Medstract, Schwartz and Hearst, BIOADI and Ab3P) for acronym extraction in Table B1 and for acronym-expansion pair extraction in Table B2. We report Glossary-level Precision, Recall, and F1-measure values for batch acronym and pair for SciAI dataset in Table B3 and for End-to-end dataset in Table B4.

¹⁶<https://forms.gle/hWJR2K64XzjpYrYY9>

¹⁷<https://forms.gle/VH1ScF2nr1PBAZK18>

Acronym and In-expansion Technique	Acronym														
	Medstract			SH			BIOADI			Ab3P			Average		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
SH	100.00%	88.63%	93.97%	99.48%	79.83%	88.58%	98.69%	78.26%	87.29%	98.20%	77.84%	86.84%	99.09%	81.14%	89.17%
MadDog	100.00%	63.63%	77.78%	95.96%	48.97%	64.85%	97.76%	54.24%	69.77%	99.13%	64.77%	78.35%	98.21%	57.90%	72.69%
SciBERT	80.64%	56.81%	66.67%	81.46%	68.72%	74.55%	85.51%	74.53%	79.64%	85.18%	71.87%	77.96%	83.20%	67.98%	74.71%
SciBERT with External Data	86.48%	72.72%	79.01%	84.47%	76.13%	80.08%	88.78%	78.67%	83.42%	87.01%	76.13%	81.21%	86.69%	75.91%	80.93%
SciDr	93.10%	61.36%	73.97%	84.30%	59.67%	69.87%	90.40%	64.38%	75.21%	88.56%	68.18%	77.04%	89.09%	63.40%	74.02%
SciDr-in with External Data	91.89%	77.27%	83.95%	91.74%	77.77%	84.18%	92.10%	79.71%	85.46%	88.11%	71.59%	78.99%	90.96%	76.59%	83.15%

Table B1: In-expansion techniques Glossary-level Precision, Recall and F1-measures for acronym for each biomedical dataset (Medstract, SH, BIOADI, and Ab3p) and the averages.

Acronym and In-expansion Technique	Pair														
	Medstract			SH			BIOADI			Ab3P			Average		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1	P	R	F1
SH	100.00%	88.63%	93.97%	95.38%	76.54%	84.93%	93.99%	74.53%	83.14%	94.98%	75.28%	83.99%	96.09%	78.75%	86.51%
MadDog	92.85%	59.09%	72.22%	91.93%	46.91%	62.12%	87.31%	48.44%	62.31%	95.21%	62.21%	75.25%	91.83%	54.16%	67.98%
SciBERT	64.51%	45.45%	53.33%	69.26%	58.43%	63.39%	66.27%	57.76%	61.72%	74.41%	62.78%	68.10%	68.61%	56.11%	61.64%
SciBERT with External Data	75.67%	63.63%	69.13%	74.42%	67.07%	70.56%	72.89%	64.59%	68.49%	76.29%	66.76%	71.21%	74.82%	65.51%	69.85%
SciDr	79.31%	52.27%	63.01%	71.51%	50.61%	59.27%	75.29%	53.62%	62.63%	80.81%	62.21%	70.30%	76.73%	54.68%	63.80%
SciDr-in with External Data	91.89%	77.27%	83.95%	81.55%	69.13%	74.83%	85.16%	73.70%	79.02%	79.72%	64.77%	71.47%	84.58%	71.22%	77.32%

Table B2: In-expansion techniques Glossary-level Precision, Recall and F1-measures for pair for each biomedical dataset (Medstract, SH, BIOADI, and Ab3p) and the averages.

Acronym and In-expansion Technique	SciAI					
	Acronym			Pair		
	P	R	F1	P	R	F1
SH	94.79%	82.63%	88.29%	91.06%	79.38%	84.82%
MadDog-in	98.43%	85.97%	91.78%	96.37%	84.17%	89.86%
SciBERT	94.63%	93.58%	94.10%	90.57%	89.56%	90.06%
SciBERT with External data	95.29%	93.49%	94.38%	91.10%	89.39%	90.24%
SciDr-in	96.75%	91.78%	94.20%	93.41%	88.62%	90.95%
SciDr-in with External data	96.91%	91.36%	94.05%	92.46%	87.16%	89.74%

Table B3: In-expansion techniques Glossary-level Precision, Recall and F1-measures for acronym and pair extraction and for the SciAI dataset.

Acronym and In-expansion Technique	User-Generated					
	Acronym			Pair		
	P	R	F1	P	R	F1
SH	90.42%	70.83%	79.43%	85.10%	66.67%	74.76%
MadDog-in	92.22%	69.17%	79.04%	87.78%	65.83%	75.23%
SciBERT	64.13%	49.16%	55.66%	56.52%	43.33%	49.05%
SciBERT with External data	48.63%	59.17%	53.38%	44.52%	54.17%	48.87%
SciDr-in	76.67%	57.50%	65.71%	67.78%	50.83%	58.09%
SciDr-in with External data	85.54%	59.17%	69.95%	80.72%	55.83%	66.01%

Table B4: In-expansion techniques Glossary-level Precision, Recall and F1-measures for acronym and pair extraction for the User Generated dataset.