

Is NLP Ready for Standardization?

Anonymous EMNLP submission

Abstract

While standardization is a well-established activity in other scientific fields such as telecommunications, networks or multimedia, in the field of AI and more specifically NLP it is still at its dawn. In this paper, we explore how various aspects of NLP (evaluation, data, tasks...) lack standards and how that can impact science, but also the society, the industry, and regulations. We argue that the numerous initiatives to rationalize the field and establish good practices are only the first step, and developing formal standards remains needed to bring further clarity to NLP research and industry, at a time where this community faces various crises regarding ethics or reproducibility. We thus encourage NLP researchers to contribute to existing and upcoming standardization projects, so that they can express their needs and concerns, while sharing their exper-

1 Introduction

Most of the Natural Language Processing community remains estranged from standardization. As this is already common practice in many computer science fields, including telecommunications, networks and multimedia, what is making NLP so special in that regard? Zielke (2020) has already asked the question of the potential barriers and benefits of standardization work in the broader field of Artificial Intelligence, which is now becoming a reality. Here we propose to deepen that discussion by investigating the more specific context of NLP and its standardization needs.

Standards are normative documents produced by Standards Developing Organizations (SDOs) such as ISO. In practice, they can be of various nature and content. Some of them are terminological references that establish shared terms and definitions for a technical domain. For instance, the ISO Online Browsing Platform¹ indexes all

¹<https://www.iso.org/obp>

existing ISO definitions. Other standards rather describe a reference framework, which can prove useful for bootstrapping new activities or rationalizing existing ones. Standards can also provide technical specifications for data, systems or procedures. This notably includes quality specifications, as well as interoperability ones (APIs, protocols, etc.). For instance, the C language (ISO/IEC 9899), the MP3 coding (ISO/IEC 11172-3), the Latin-1 charset (ISO/IEC 8859-1) and the 2-letter language codes (ISO 639-1) are all examples of standards.

Standards are written by volunteer experts from various backgrounds (scientific, legal, standardization experts, etc.). In most SDOs, registration is open to anyone willing to contribute, usually through a mirror committee within their national standards organization. Experts collaborate within working groups and decisions are taken by consensus² – across countries, but also across backgrounds, across sectors, across technical fields. This approach makes standardization a rather slow process (with up to 3 years to establish some standards), but it also ensures the strength of the agreements.

Standards are especially important for the industry and for regulatory authorities – but they can apply on very technical fields, including scientific topics, and therefore affect the research community as well. This paper investigates how NLP standardization could impact our community, offering both challenges and opportunities.

After a brief review of the current state of NLP standardization initiatives (§2), we explore standardization gaps pertaining to NLP evaluation (§3), data and formats (§4), tasks (§5) and higher-level

²Compared to unanimity, where everyone supports the decision, consensus means that noone objects to the decision. This decision process helps to identify middle ground solutions that everyone in an heterogeneous group can find acceptable, whereas an unanimity requirement would bear the risk of freezing projects due to unsolvable cultural differences or diverging interests.

076 concepts (§6). Having built a broader view of how
077 NLP standardization could benefit research, but
078 also the society, industry and regulations (§7), we
079 then conclude on possible contributions that NLP
080 researchers could add to those efforts (§8).

081 **2 Existing initiatives towards NLP** 082 **standardization**

083 Within the NLP community, most of the stan-
084 dardized material is actually *de facto* standards:
085 data, tools or methodology that are consensually
086 used throughout the field, even though they don't
087 have any official status and have not necessarily
088 gone through the formalization process that offi-
089 cial standards offer. Such *de facto* standards of-
090 ten result from past shared tasks: for instance,
091 the `mteval-v13a.pl` evaluation script from the
092 WMT shared tasks series has been used for years
093 as the reference evaluation script by a large part
094 of machine translation research. Similarly, the
095 CoNLL-X format for dependency treebanks has
096 been widely adopted following the corresponding
097 CoNLL shared task (Buchholz and Marsi, 2006).
098 As for event detection, the definition of the task
099 itself is fully driven by the ACE 2005 campaign
100 (Walker et al., 2006), to the point that it is some-
101 times referred to as “ACE event detection” (Chen
102 et al., 2018). Recent years have seen however the
103 growth of the Universal Dependencies initiative
104 (Nivre et al., 2016), for establishing common guide-
105 lines for treebank annotation across languages;
106 considering the breadth of its contributors and its
107 sustained efforts for guidelines formalization, this
108 project has now become very close to standardiza-
109 tion work and could be considered as an SDO.

110 Regarding official standardization initiatives rel-
111 evant for NLP, the most established one is ISO's
112 Technical Committee 37 (*Language and Terminol-
113 ogy*), and more prominently its subcommittees 4
114 (*Language resource management*, created in 2001)
115 and 5 (*Translation, interpreting and related tech-
116 nology*, created in 2012). With a strong focus on
117 corpora and annotation, these groups have notably
118 released a number of annotation framework stan-
119 dards (e.g. for TIGER-XML or TEI), which are
120 extensively used by the corresponding industry;
121 they also co-organize with ACL the ISA workshop
122 series on Interoperable Semantic Annotation (Bunt,
123 2021). Yet this focus on data leaves the algorithmic
124 and evaluation parts of NLP largely unaddressed.

125 ISO-IEC's Joint Technical Committee 1 (*Infor-*

126 *mation Technology*) has created in 2017 its sub-
127 committee 42 on Artificial Intelligence. This one
128 is more concerned with algorithms, development
129 methodology, and system evaluation, but at the
130 higher level of AI in general and not delving into
131 NLP-specific aspects. To date, its NLP-related ac-
132 tivity has focused on defining a few major concepts,
133 such as NLG, question answering, or OCR. Concur-
134 rently, other global SDOs such as ITU-T have also
135 explored NLP standardization, but mostly in the
136 context of specific use cases (e.g. ITU-T H.862.5
137 “*Emotion enabled multimodal user interface based
138 on artificial neural networks*”) rather than the NLP
139 field in general. Hence their remains a gap in terms
140 of NLP standardization.

141 The European counterparts of ISO-IEC (CEN-
142 CENELEC) have thus created in 2021 their own
143 Joint Technical Committee 21 on Artificial Intel-
144 ligence, with a group dedicated to kickstarting ac-
145 tivities on speech and NLP (ad-hoc group 4, *AI
146 systems for human language processing*, on track
147 to become persistent in 2022). This is the first page
148 of a new chapter, roadmaps are being written today.
149 Now is thus the right time to question what can,
150 and should, be standardized within our field.

151 **3 Is NLP evaluation ready for** 152 **standardization?**

153 The reproducibility crisis that has spread through
154 the field in recent years (Belz et al., 2021b; Lu-
155 cic et al., 2022) has renewed the community's in-
156 terest for fair and reliable evaluation. This has
157 led in the 2020s to a blooming of workshops
158 and shared tasks dedicated to evaluation means
159 (whether human evaluation or automated metrics),
160 such as Eval4NLP, HumEval, GEM, or *Benchmark-
161 ing: Past, Present and Future* (Eger et al., 2020;
162 Belz et al., 2021a; Bosselut et al., 2021; Church
163 et al., 2021). But those concerns are not new, as
164 illustrated by the 4REAL workshops organized in
165 the 2010s (Branco et al., 2016). Even the terminol-
166 ogy of reproducibility has been the topic of debate
167 and clarification attempts for many years in the
168 machine learning community (Drummond, 2009). The
169 existence of the LREC conference series itself is a
170 token of that interest for standardized evaluation,
171 with its first edition dedicating a whole workshop
172 to the lack of a shared strategy, definition and in-
173 frastructure for system evaluation (ELSE, 1998;
174 McTait and Choukri, 2003).

175 For domains like human evaluation of NLG, the

need for more consistent practices is clear enough, and [Howcroft et al. \(2020\)](#) have already advocated for producing standards on both the methodology and the terminology, based on their review of 20 years of NLG with conflicting evaluation criteria. Yet here we argue that even cases with seemingly straightforward automated evaluation can suffer today from the lack of standards.

3.1 On defining metrics

One of the challenges of standardized evaluation is to ensure that the metrics used are defined in a way that leaves no place to ambiguity, which in practice is rarely the case in the field. For instance, even the well-known F1-score, despite its very formal definition as the harmonic mean of precision ($\frac{TP}{TP+FP}$) and recall ($\frac{TP}{TP+FN}$), becomes ambiguous when applied to tasks such as Named Entity Recognition.

A first issue resides in the common practice of casting this chunking task into a sequence labelling task, through BIO-style token-level encoding: B-... labels denote the first token of an entity, I-... labels other tokens within the entity, and O labels other tokens. This raises the question of which objects are considered for true/false positives/negatives: those labels, or the chunks. For instance, with B-PER I-PER O B-LOC O as a reference sequence, predicting B-PER O O B-LOC I-LOC yields a score of 60 (micro-)F1 if evaluated as a sequence labeling task (67 F1 if O is not considered a class) by looking at tokens, but 0 F1 if evaluated as a chunking task by looking at the predicted chunks (here with exact match). While most experienced researchers know to prefer the latter (and know where to get that information), the youngest researchers as well as industry practitioners are not necessarily aware of this implicit rule. Such confusion can in turn lead to incorrect comparison of models, or incorrect reporting of product performance.

[Taillé et al. \(2020\)](#) report other underspecified aspects, such as the criteria to accept true positives (with or without typing, with partial or exact match...), the use of micro- or macro-averaging, or the existing practice to ignore some classes (such as Other or MISC). As they highlight, these issues also propagate to evaluation of relation extraction, and just one of those can already lead to overestimating the results by up to +3 F1 on a widely-used dataset.

3.2 On implementing metrics

Another challenge is the underspecification of implementation details for those evaluation metrics, even when the metric itself has a non-ambiguous definition.

For machine translation, [Post \(2018\)](#) investigates the divergence in scores that can result from different implementations of the BLEU metric, based on diverging choices of parameters and preprocessing (e.g. the maximum n-gram length, the number of references, or user-supplied and/or metric-internal tokenization). He reports up to 1.8 BLEU difference when varying only the tokenization used for scoring, which is actually more than the gains measured for BPE ([Sennrich et al., 2016](#)), which was a game changer for neural machine translation.

Such variations in implementation can occur even in cases as seemingly simple as using F1-scores for classification: for instance [Belz \(2021\)](#) compare concurrent reproduction studies of the same text classifier, and report score divergences up to 5.2 F1 due to metric reimplementations.

Another source of implementation divergence is the procedure adopted to deal with invalid outputs (ill-formatted, impossible sequences, etc.). In the case of Named Entity Recognition, [Lignos and Kamyab \(2020\)](#) investigate how different strategies to repair invalid BIO sequences within the scorer can impact the measured F1, a condition which according to [Palen-Michel et al. \(2021\)](#) also affects the gold labels in a number of renowned datasets, and leads to differences up to 3.25 F1 in a realistic scenario. For the BIOES encoding scheme alone, [Kroutikov \(2019\)](#) numbers at least 7776 different strategies that could be adopted to repair invalid label pairs.

3.3 Tooling to the rescue?

As a means to circumvent those pitfalls, [Lignos and Kamyab \(2020\)](#) advocates for never reimplementing evaluation metrics and relying instead on third-party reference tools. This is in line with [Post \(2018\)](#)'s strategy to release the SacreBLEU package, with the hope that its configurability, documentation, ease of use and variant reporting will enable standardized evaluation. Can tools alone indeed fill in for standards?

The main issue with that view is that it supposes that tools are correctly used. However, [Marie et al. \(2021\)](#) unveil that the growing number of users of SacreBLEU are in practice often misusing it (not re-

275 porting the variant used, comparing its scores with
276 other scorers, etc.). Similarly, Palen-Michel et al.
277 (2021) release SeqScore as a possible reference
278 tool for named entity recognition evaluation, but
279 they do so based on the failure of previous *de facto*
280 standard tools. For instance, Akbik et al. (2019) ob-
281 serve that their previous paper (which has now over
282 1000 citations) had overestimated its results by up
283 to 0.8 F1, because they used the official CoNLL-03
284 evaluation script (designed for BIO) on a BIOES-
285 encoded dataset. On a side note, it can also happen
286 that the most popular scorer simply contains a bug
287 – how can this be assessed if the tool itself serves
288 as standard?

289 Another possible approach would be to rely
290 more heavily on Kaggle-style benchmarking plat-
291 forms that enable fairer comparison than stan-
292 dalone evaluation tools. The issue here is that such
293 practices can arbitrarily foster inadequate evalua-
294 tion. Bowman and Dahl (2021) now consider NLU
295 evaluation “broken” due to benchmark-driven stan-
296 dardization of practices: a number of those bench-
297 marks are actually rewarding unreliable and biased
298 systems. They leave no place to reflect upon a
299 given system’s appropriate evaluation setting, and
300 instead incentivize gaming the numbers. Church
301 and Hestness (2019) review 25 years of evaluation
302 practices and show how the rigour efforts that have
303 led to such benchmarks are now pushing against
304 their initial purpose of bringing more insights to
305 “content-free debates”. Overall, leaderboards have
306 drawn a lot of criticism in recent years (Rogers,
307 2019; Ethayarajh and Jurafsky, 2020; Kiela et al.,
308 2021) and are therefore a poor candidate to address
309 the lack of standards.

310 Instead of producing and relying on tools, typical
311 standardization work would rather approach the
312 issue by writing comprehensive specifications of
313 the evaluation metrics (detailing their computation,
314 their usage, their meaning), which can in turn apply
315 on tools. This includes providing the means to
316 verify that a given scorer or a given evaluation
317 protocol is compliant with the specification. Hence
318 comparable evaluation can be formally ensured, but
319 not at the cost of insights and appropriateness.

320 3.4 Does it matter?

321 So the lack of standards leads to more imprecision
322 in the measures and less rigorous comparisons. Is
323 that really an issue, as long as those numbers are
324 high and continue increasing, whatever the criteria?

325 Haven’t experimental sciences handled imprecision
326 for centuries, and accepted that challenge as part
327 of the job?

328 According to Morey et al. (2017), such imprec-
329 sion has already endangered scientific progress
330 in whole fields of NLP: in their review of several
331 years of contributions in discourse parsing, they
332 discover that the various conclusions drawn on the
333 benefits of distributed representations are mostly
334 wrong in that field. What was considered a huge
335 improvement, with 24 to 51% relative error reduc-
336 tion depending on the metric, was actually a gain
337 of 11 and 16% for two of the metrics, and a *loss*
338 of 15 and 53% for the other two. Here the culprit
339 was the choice to macro-average over documents in
340 some but not all of the works, following practices
341 existing in different communities.

342 The lack of standards can thus lead to misinter-
343 preting regress as progress. But it can also affect
344 the wider world outside of research. For instance
345 when a contract is signed, and B2B products are
346 to be developed according to a given performance
347 level specified in the contract, there should be no
348 place to ambiguity. Who should be the judge of
349 whether the contract is fulfilled, if the bar is met by
350 one implementation variant of the metric, but not
351 the other? And what if a regulation contains such
352 performance requirement?

353 Comparability is also a strong enabler for indi-
354 vidual rights as consumers. Potential users should
355 be able to make an informed choice when com-
356 paring existing products. Transparency regulations
357 can contribute to that, but that information becomes
358 meaningless if the same number can be interpreted
359 differently depending on implementation details.

360 3.5 Can standards hinder research?

361 Scientific concerns regarding NLP evaluation go
362 in fact way beyond the need for fair comparison.
363 A number of automated metrics in wide use to-
364 day have poor correlation with human judgment,
365 and a lot of research efforts have been devoted
366 to designing more relevant metrics. Notoriously,
367 WMT has been running an annual shared task on
368 machine translation evaluation metrics since 2008
369 (Bojar et al., 2016; Mathur et al., 2020; Freitag
370 et al., 2021), thereby consolidating the community
371 consensus that the BLEU metric certainly has its
372 utility, but also a number of shortcomings (Reiter,
373 2018), and it is far from being the best metric in
374 existence. METEOR, chrF, CharacTer, BERTscore

(Banerjee and Lavie, 2005; Popović, 2015; Wang et al., 2016; Zhang et al., 2019) are just a few examples among a broad panel of often more appropriate metrics, even though the single-best “one-size-fits-all” metric has not been found yet.

One possible fear with standardization could then be to prevent researchers from pursuing their quest for the best metric, or simply to prevent them from using in their work another good metric instead of BLEU – leading again to fostering bad evaluation practices and limiting the insights brought to future research. However, standards do not need to be compulsory. It is quite possible to write them in a way that preserves that research freedom, but still brings some order and clarity. BLEU is not the best metric, but BLEU is nevertheless preferable to exotic approaches such as measuring an F1-score at the sentence level (true positive if the sentence is an exact match). Are we confident that all practitioners that may have to evaluate a machine translation system at some point (including e.g. software developers in the industry) are aware of that? Can we at least give formal existence to that tiny piece of knowledge?

It is indeed a fact that in a number of cases in NLP evaluation, it is not necessarily known what is the most appropriate choice among the various existing variants. It can also be use case-dependent. And standards in such context are not meant to arbitrarily foster one option among the others. Their role here would rather be to formally reference and specify the existing relevant options (pushing away the ones that are already known to be inappropriate), and offer practical ways to declare, identify or verify which one of those options has indeed been used in a given paper or product.

4 Are NLP data and formats ready for standardization?

Yes they are, and they have been as early as 1993, when EAGLES (Expert Advisory Group for Language Engineering Standards) was established to develop such standards. Ide et al. (2017) review 30 years of community progress from confusion to *de facto* standards to standards. However, despite marked efforts from ISO’s Technical Committee 37, this paradigm has only been adopted so far in some parts of the field, and much progress remains to achieve for fully standardizing NLP annotations.

In particular, Ide et al. (2017) underline the need to better standardize the *content* of annotations.

While many (although not all) corpus authors have gone through the formalization process of writing annotation guidelines, this has mostly led to a profusion of co-existing guidelines for the same task. The case of dependency parsing is interesting in that regard, as the Universal Dependencies project managed to unify most of the pre-existing annotation schemes, while preserving their idiosyncrasies (Nivre et al., 2016). Yet this is a success story that most parts of the field have not had so far.

In addition, annotation processes should include some quality control mechanisms, such as measuring inter-annotator agreement (Hovy and Lavid, 2010). However, there is poor consensus on what would be a “good” agreement value for a given task, depending on its complexity and subjectivity (Artstein and Poesio, 2008; Mathet et al., 2012). Are we even sure that inter-annotator agreement is an appropriate quality control (Wong and Lee, 2013; Passonneau and Carpenter, 2014; Plank et al., 2014; Boguslav and Cohen, 2017; Basile et al., 2021)? In recent years, the growing reliance on crowdsourcing has only strengthened the challenges, hence the pressing need for standardizing practices (Sabou et al., 2014).

Standardization gaps do not concern solely the semantics of the annotation, but also their format. Looking at machine translation, parallel corpus formats include SGML (for which WMT maintains a `wrap-xml.perl` script to preserve compatibility with scoring scripts), XML (with XCES for sentence alignment), TMX, bitext (two files with corresponding line numbers), but also tabular formats with per-language columns separated by either tabs or other separators. The OpusTools converters (Aulamo et al., 2020) support only part of that spectrum. As for named entity recognition, co-existing encoding schemes include IO, IOB (aka IOB1), BIO (aka IOB2), BIOES (aka IOBES), BILOU (aka BILUO) and BMEOW (Palen-Michel et al., 2021). And there are others (Malik and Sarwar, 2016). One can always write converters, but this is tedious work, and prone to introducing discrepancies in case of invalid sequences (see §3.2). Third-party open source converters can help (Lester, 2020), yet they usually support only some of the encoding schemes. Formats can further differ when considering the file format: whereas CoNLL-2003 was distributed as tabular IOB (Tjong Kim Sang and De Meulder, 2003), spaCy relies on JSON BILUO. And this is only for

sequence tag schemes, while MUC-6 uses SGML (Grishman and Sundheim, 1996) and WiNER-fr prefers an offset-based scheme to directly encode the spans (Dupont, 2019).

In terms of input and output formats, NLP tools can already rely on a number of extensible pipelines such as Stanford CoreNLP or spaCy (Manning et al., 2014; Honnibal and Montani, 2017), as well as abstraction frameworks such as AllenNLP or PyText (Gardner et al., 2018; Aly et al., 2018) – but this differs from actual APIs designed for interoperability among products. Today such interoperability is mostly fostered by infrastructure-based initiatives such as the Language Application Grid (Ide et al., 2016, 2015). The European Language Grid project (Rehm et al., 2020a, 2021) now proposes to build an umbrella platform that hosts resources but also unifies NLP APIs through its “functional services” infrastructure. In addition, Kim et al. (2020) propose to standardize a web protocol for NLPaaS, while Rehm et al. (2020b) set a roadmap of interoperability levels to enable cross-platform workflows. Instead of duplicating those projects, the role of SDOs here would rather be to build upon those APIs, by escalating them into official standards with formal specifications.

Finally, data warrants data documentation. This is another area where individual initiatives have produced valuable guidelines on necessary metadata (Bender and Friedman, 2018). But work still remains to give that material more formalism and ensure consensus across communities.

5 Are NLP tasks ready for standardization?

Getting to the core of NLP, even the tasks themselves warrant further consideration for standardization. Indeed, NLP research has recently gained awareness that making further progress on NLU tasks now meant taking some detours to better define terms like “meaning” (Bender and Koller, 2020), “comprehension” (Dunietz et al., 2020) and the associated tasks. Yet even the basic expectations on inputs/outputs can be underspecified for some tasks. For instance, question answering can refer to various concrete tasks, such as multiple-choice answer selection (Aydin et al., 2014), span extraction (with or without paragraph retrieval) in the SQuAD style (Rajpurkar et al., 2016), free-form answering that can include multi-hop ques-

tions (Chen et al., 2019), or answering questions over knowledge bases (Fu et al., 2020), which don’t warrant the same algorithmic approaches. Gardner et al. (2019) propose to solve the conundrum by considering question answering as a format and splitting it from the definition of the task; yet even then the taxonomy remains dense (Rogers et al., 2021).

Information extraction is another field where tasks and their terminology are largely ill-defined. Even its primary task, named entity recognition, has been subject to a number of conceptual debates (Marrero et al., 2013). Entity linking is better delineated, but has been associated with a number of different names: entity linking, named entity linking, named entity disambiguation, named entity normalization... Are all of those terms synonymous, or do they slightly differ in scope? The literature has already proposed many definitions for entity linking, often inconsistently: for instance Shen et al. (2014) write both “Entity linking is the task to link entity mentions in text with their corresponding entities in a knowledge base” and “to link named entity mentions appearing in web text with their corresponding entities in a knowledge base, which is called entity linking”. This notably raises doubts as to whether entity linking applies only to named entities, or also to non-named entities (Paris and Suchanek, 2021). Or to non-named mentions of entities that have names? In the lack of terminological standards, presumably the best definition of the task is to look at how the corpus at hand has been annotated; but then the task definition can vary a lot from one dataset to another, so that evaluating an entity linking approach on multiple datasets may not make actual sense. Many other discrepancies could be listed here (e.g. relation extraction referring to either relation clustering or open information extraction), but in the end, the name “information extraction” itself is an ill-defined term, with a functional scope that varies a lot depending on individuals. So if a system is branded as an information extraction system, what are its functionalities supposed to be?

Time is not innocent in those terminological conflicts. Language modelling is one striking example of terminological drift. Historically, language models meant “a probability distribution over all possible word strings in a language” (Arisoy et al., 2012) – or even a next-word predictor, as in the n-gram paradigm: “language modeling, the problem

577 of predicting the next word based on words already
578 seen before” (Xu and Jelinek, 2004). But since
579 2018 and the advent of *masked* language models,
580 the term “language model” has now shifted to refer
581 to Transformer-based contextualized embeddings,
582 regardless of any probability distribution, and not
583 necessarily autoregressive (as in [Ettinger, 2020](#)).

584 Are these discrepancies an issue? Semantic drift
585 is a natural phenomenon in any language, and a
586 profusion of definitions also means a profusion of
587 problems addressed by the community as a whole.
588 However, trouble arises when using those task defi-
589 nitions to catalog or to assess existing systems: how
590 to decide whether a given system meet one’s expect-
591 ations, if it is branded with ambiguous functional-
592 ities? Achieving clarity on product capabilities is a
593 matter of commercial interest for companies, and
594 of consumer rights for individuals. But it can also
595 affect scientific processes, as exemplified by the
596 Great Misalignment Problem ([Hämäläinen and Al-
597 najjar, 2021](#)) between blurry objectives, the actual
598 task fulfilled by the system, and the task against
599 which human evaluation is performed.

600 **6 Are NLP concepts ready for** 601 **standardization?**

602 At a higher level, a number of concepts would also
603 benefit from formal standards. This notably con-
604 cerns the term “multilingual”, which has been used
605 to describe very different properties, such as: a sys-
606 tem that juxtaposes models for multiple languages
607 ([Otero and González, 2012](#)), with or without inter-
608 nal language identification; an algorithm that does
609 not rely on language-specific features or knowl-
610 edge, and can therefore be trained on a dataset
611 from any language ([Johansson and Nugues, 2006](#);
612 [Szarvas et al., 2006](#)), even though this does not
613 guarantee actual language independence ([Bender,
614 2011](#)); or a single model that can indiscriminately
615 process contents from many languages ([Pires et al.,
616 2019](#)). Focusing only on the latter definition, how
617 many is many? And how diverse? Can an Indo-
618 European-only system be considered multilingual?
619 In light of rising initiatives for fostering more lan-
620 guage diversity in NLP research ([Bender, 2019](#);
621 [Joshi et al., 2020](#)), including a dedicated theme
622 track at ACL 2022, it now appears pressing to es-
623 tablish consensual criteria on what renders a given
624 system multilingual. Otherwise, how can progress
625 in that matter be quantified?

626 Trustworthiness is another relevant concept for

627 NLP systems, especially from the viewpoint of
628 policy makers. The [High-Level Expert Group on
629 AI \(2019\)](#) has notoriously established a list of AI
630 trustworthiness characteristics, but they still lack
631 shared actionable definitions. The concept of bias
632 for instance, while subject to a growing interest
633 in NLP research, is rarely formally defined in that
634 literature, or with diverging senses ([Blodgett et al.,
635 2020](#)), even though that conceptualization should
636 be a prerequisite before defining the corresponding
637 bias measures ([Dev et al., 2021](#)). “Robustness”
638 is similarly overloaded, with meanings ranging
639 from maintained performance on out-of-domain
640 data ([Bernier-Colborne and Langlais, 2020](#)), on
641 transformed data ([Sanchez et al., 2018](#); [Gan and
642 Ng, 2019](#)), or in presence of natural noise ([Zhou
643 et al., 2019](#)), to specific defenses against adversarial
644 attacks ([Hsieh et al., 2019](#)). A fortiori, there is no
645 formal taxonomy on what kind of noise a “robust”
646 NLP system should minimally handle: typos only,
647 or L2 learners grammar errors, lexical borrowings?
648 Broken encoding? Or others?

649 Regarding interpretability and explainability,
650 while some use those terms interchangeably, oth-
651 ers have drawn firm distinctions: interpretability is
652 “loosely defined as the science of comprehending
653 what a model did (or might have done)” ([Gilpin
654 et al., 2018](#)), while “Given a certain audience,
655 explainability refers to the details and reasons a
656 model gives to make its functioning clear or easy to
657 understand” ([Arrieta et al., 2020](#)), thereby putting
658 the cognitive load of that understanding process
659 more on the model and less on the human. As it
660 seems that “interpretability” has become the pre-
661 ferred term in the NLP community, while other
662 fields rather use “explainability”, should that dif-
663 ference of focus be understood as a conceptual
664 divergence of interests (whereby the NLP commu-
665 nity would foster more involvement of the human
666 in model understanding than other communities),
667 or only as a terminological discrepancy? Concur-
668 rently, “explainability” as expressed by some other
669 audiences (especially non-technical ones) has noth-
670 ing to do with either of those concepts, and is rather
671 a synonym for “transparency”, “testing”, or even
672 “reproducibility” ([Brennen, 2020](#)), hence a looming
673 crisis if policy makers mean one while practitioners
674 understand the other. As for transparency, while
675 there is consensus on the need for auditability and
676 documentation, the question of what has to be doc-
677 umented and how is still open ([Saxon et al., 2021](#)).

678 **7 On the benefits of NLP standardization**

679 Looking at the long-term impact of the Universal
680 Dependencies project hints at how standardizing
681 NLP could more generally benefit the field and its
682 dynamics. The immediate benefit was a more faith-
683 ful evaluation across languages, that enabled deeper
684 investigation of cross-lingual transfer. But the exist-
685 ence of common guidelines also created a commu-
686 nity incentive to produce more data, by providing
687 the guidance and means for easier extension to
688 dozens of zero-resourced languages. The creation
689 of the project itself opened new fora for commu-
690 nity-level collaboration and sharing, thereby giving de-
691 pendency research a new boost. This has been
692 an opportunity to reflect upon research practices,
693 fostering more systematic studies on annotation
694 scheme impact, better highlighting the gaps in lin-
695 guistic coverage, and uncovering biases in our view
696 of syntax. Overall, standardization contributes to
697 better driving research, both at the individual and
698 institutional levels. A clear taxonomy makes it
699 easier to identify scientific gaps, more compati-
700 ble resources and tools offer richer experimental
701 means, and shared definitions guarantee that we are
702 all pushing in the same direction.

703 Standards also support community-wide adop-
704 tion of good practices. Even if abundantly dis-
705 cussed and well documented as checklists, espe-
706 cially regarding evaluation (Ribeiro et al., 2020;
707 van der Lee et al., 2019; Gehrmann et al., 2022;
708 Marie et al., 2021; Escartín et al., 2021) and doc-
709 umentation (Bender and Friedman, 2018; Gebru
710 et al., 2021; Mitchell et al., 2019; Ligozat and Luc-
711 cioni, 2021; Wilkinson et al., 2016), consensual
712 good practices guidelines are not necessarily im-
713 plemented in practice, in research and even more
714 crucially in industry. Escalating them to formal
715 standards makes it easier to enforce them.

716 Comparability is another clear benefit for NLP
717 researchers, but even more so for users and con-
718 sumers. Interoperability can facilitate putting a
719 researcher’s ideas into users’ hands, with easier
720 integration into products. But it is also a matter
721 of survival for SMEs, for packaging and distribut-
722 ing their products in a competitive environment
723 where Big Tech standalone solutions dominate the
724 market and SMEs struggle to propose large-scale
725 alternatives.

726 Last but not least, standardizing NLP concepts
727 is a necessary step to refine a shared roadmap to
728 address ethical considerations in NLP (Hovy and

Spruit, 2016; Leidner and Plachouras, 2017) – but
it is also a prerequisite to regulating abuses and
enforcing safe use of NLP that preserves individ-
ual rights. At a time where the EU is establishing
its AI Act, it has started collaborating with CEN-
CENELEC to bring clarity on the terminology and
processes needed to legally ensure key trust char-
acteristics, such as robustness, transparency or fair-
ness.

8 Last words: are we ready?

This review of various aspects of the NLP ecosys-
tem has shown how the lack of standards can cause
confusion, inefficiency, and sometimes even render
research efforts detrimental to scientific progress,
through misinterpretations and fostering of bad
practices. Building and formalizing consensus on
key practices and concepts would enable instead a
more reproducible, more insightful, more industry-
ready and more ethical science.

Admittedly, not everything can be readily stan-
dardized: sometimes the scientific material to do
so does not even exist yet. And sometimes research
freedom and creativity need to be preserved by
maintaining concurrent options. But these are pre-
cisely cases where it is even more important that
the standardization ecosystem benefits from scien-
tific expertise, in order to avoid over-standardizing
the field, or widening the discrepancies between
research and industry practices.

We believe it is a matter of scientific responsibil-
ity to offer such guidance to those who are shaping
the industrial and legal future of society-wide use
of NLP. Contributing to standardization means shar-
ing our expertise and insights, but also our needs
and our concerns, both as scientists and as citizens.
NLP is ready and in need – now *we* have to get
ready.

There are numerous ways to taking part. While
community-internal initiatives should be pursued
and fostered, we also encourage European re-
searchers to join CEN-CENELEC/JTC 21 for con-
tributing to its budding roadmap, and worldwide
researchers to both pursue resource standardization
efforts within ISO/TC 37 and help ISO-IEC/JTC
1/SC 42 to deepen debates that have still only
scratched the surface of the upcoming work. Or-
ganize events, discuss, share, debate, draft, brain-
storm, publish. And NLP standards will be within
reach.

778 Limitations

779 A significant part of this paper has a purely illustrative value, and the provided set of examples does
780 not convey a comprehensive view of the existing
781 standardization issues. Similarly, despite extensive
782 search, we offer no guarantee of exhaustivity in
783 our inventory of NLP standardization groups, in
784 particular for non-cited SDOs (e.g. IEEE).
785

786 The review and discussion are also biased to-
787 wards a number of European concerns and initia-
788 tives, which may be either a symptom of its pio-
789 neering position on the topic, or merely a lack of
790 depth in our survey of local initiatives in other parts
791 of the world. National-level standardization efforts
792 are not discussed either.

793 Finally, this work only scratches the surface of
794 discussing the scientific and industrial feasibility
795 of standardization for each part of the field, which
796 may significantly vary from one task or concept to
797 another, depending on their maturity and history.

798 References

799 Alan Akbik, Tanja Bergmann, and Roland Vollgraf.
800 2019. [Pooled contextualized embeddings for named
801 entity recognition](#). In *Proceedings of the 2019 Con-
802 ference of the North American Chapter of the Asso-
803 ciation for Computational Linguistics: Human Lan-
804 guage Technologies, Volume 1 (Long and Short Pa-
805 pers)*, pages 724–728, Minneapolis, Minnesota. As-
806 sociation for Computational Linguistics.

807 Ahmed Aly, Kushal Lakhota, Shicong Zhao, Mri-
808 nal Mohit, Barlas Oguz, Abhinav Arora, Sonal
809 Gupta, Christopher Dewan, Stef Nelson-Lindall, and
810 Rushin Shah. 2018. [Pytext: A seamless path
811 from nlp research to production](#). *arXiv preprint
812 arXiv:1812.08729*.

813 Ebru Arisoy, Tara N. Sainath, Brian Kingsbury, and
814 Bhuvana Ramabhadran. 2012. [Deep neural network
815 language models](#). In *Proceedings of the NAACL-
816 HLT 2012 Workshop: Will We Ever Really Replace
817 the N-gram Model? On the Future of Language
818 Modeling for HLT*, pages 20–28, Montréal, Canada.
819 Association for Computational Linguistics.

820 Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez,
821 Javier Del Ser, Adrien Bannetot, Siham Tabik, Al-
822 berto Barbado, Salvador García, Sergio Gil-López,
823 Daniel Molina, Richard Benjamins, et al. 2020. [Ex-
824 plainable artificial intelligence \(xai\): Concepts, tax-
825 onomies, opportunities and challenges toward re-
826 sponsible ai](#). *Information fusion*, 58:82–115.

827 Ron Artstein and Massimo Poesio. 2008. [Survey ar-
828 ticle: Inter-coder agreement for computational lin-
829 guistics](#). *Computational Linguistics*, 34(4):555–
830 596.

Mikko Aulamo, Umut Sulubacak, Sami Virpioja, and
Jörg Tiedemann. 2020. [OpusTools and parallel cor-
pus diagnostics](#). In *Proceedings of the 12th Lan-
guage Resources and Evaluation Conference*, pages
3782–3789, Marseille, France. European Language
Resources Association.

Bahadir Ismail Aydin, Yavuz Selim Yilmaz, Yaliang Li,
Qi Li, Jing Gao, and Murat Demirbas. 2014. [Crowd-
sourcing for multiple-choice question answering](#). In
AAAI, pages 2946–2953. Citeseer.

Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:
An automatic metric for MT evaluation with im-
proved correlation with human judgments](#). In *Pro-
ceedings of the ACL Workshop on Intrinsic and Ex-
trinsic Evaluation Measures for Machine Transla-
tion and/or Summarization*, pages 65–72, Ann Ar-
bor, Michigan. Association for Computational Lin-
guistics.

Valerio Basile, Michael Fell, Tommaso Fornaciari,
Dirk Hovy, Silviu Paun, Barbara Plank, Massimo
Poesio, and Alexandra Uma. 2021. [We need to con-
sider disagreement in evaluation](#). In *Proceedings of
the 1st Workshop on Benchmarking: Past, Present
and Future*, pages 15–21, Online. Association for
Computational Linguistics.

Anya Belz. 2021. [Quantifying reproducibility in nlp
and ml](#). *arXiv preprint arXiv:2109.01211*.

Anya Belz, Shubham Agarwal, Yvette Graham, Ehud
Reiter, and Anastasia Shimorina, editors. 2021a. [Proceedings of the Workshop on Human Evaluation
of NLP Systems \(HumEval\)](#). Association for Compu-
tational Linguistics, Online.

Anya Belz, Shubham Agarwal, Anastasia Shimorina,
and Ehud Reiter. 2021b. [A systematic review of re-
producibility research in natural language process-
ing](#). In *Proceedings of the 16th Conference of the
European Chapter of the Association for Computa-
tional Linguistics: Main Volume*, pages 381–393,
Online. Association for Computational Linguistics.

Emily Bender. 2019. [The# benderrule: On naming the
languages we study and why it matters](#). *The Gradi-
ent*, 14.

Emily M Bender. 2011. [On achieving and evaluating
language-independence in nlp](#). *Linguistic Issues in
Language Technology*, 6.

Emily M. Bender and Batya Friedman. 2018. [Data
statements for natural language processing: Toward
mitigating system bias and enabling better science](#).
*Transactions of the Association for Computational
Linguistics*, 6:587–604.

Emily M. Bender and Alexander Koller. 2020. [Climb-
ing towards NLU: On meaning, form, and under-
standing in the age of data](#). In *Proceedings of the
58th Annual Meeting of the Association for Compu-
tational Linguistics*, pages 5185–5198, Online. As-
sociation for Computational Linguistics.

887	Gabriel Bernier-Colborne and Phillippe Langlais. 2020.	Harry Bunt, editor. 2021. <i>Proceedings of the 17th Joint</i>	943
888	HardEval: Focusing on challenging tokens to assess	<i>ACL - ISO Workshop on Interoperable Semantic An-</i>	944
889	robustness of NER . In <i>Proceedings of the 12th Lan-</i>	<i>notation</i> . Association for Computational Linguistics,	945
890	<i>guage Resources and Evaluation Conference</i> , pages	Groningen, The Netherlands (online).	946
891	1704–1711, Marseille, France. European Language		
892	Resources Association.		
893	Su Lin Blodgett, Solon Barocas, Hal Daumé III, and	Anthony Chen, Gabriel Stanovsky, Sameer Singh, and	947
894	Hanna Wallach. 2020. Language (technology) is	Matt Gardner. 2019. Evaluating question answer-	948
895	power: A critical survey of “bias” in NLP . In <i>Pro-</i>	ing evaluation . In <i>Proceedings of the 2nd Workshop</i>	949
896	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	<i>on Machine Reading for Question Answering</i> , pages	950
897	<i>ciation for Computational Linguistics</i> , pages 5454–	119–124, Hong Kong, China. Association for Com-	951
898	5476, Online. Association for Computational Lin-	putational Linguistics.	952
899	guistics.		
900	Mayla Boguslav and Kevin Bretonnel Cohen. 2017.	Yubo Chen, Hang Yang, Kang Liu, Jun Zhao, and	953
901	Inter-annotator agreement and the upper limit on ma-	Yantao Jia. 2018. Collective event detection via a	954
902	chine performance: Evidence from biomedical natu-	hierarchical and bias tagging networks with gated	955
903	ral language processing. In <i>MEDINFO 2017: Pre-</i>	multi-level attention mechanisms . In <i>Proceedings of</i>	956
904	<i>cision Healthcare through Informatics</i> , pages 298–	<i>the 2018 Conference on Empirical Methods in Nat-</i>	957
905	302. IOS Press.	<i>ural Language Processing</i> , pages 1267–1276, Brus-	958
906	Ondrej Bojar, Christian Federmann, Barry Haddow,	sels, Belgium. Association for Computational Lin-	959
907	Philipp Koehn, Matt Post, and Lucia Specia. 2016.	guistics.	960
908	Ten years of wmt evaluation campaigns: Lessons	Kenneth Church, Mark Liberman, and Valia Kordoni,	961
909	learnt. In <i>Proceedings of the LREC 2016 Work-</i>	editors. 2021. <i>Proceedings of the 1st Workshop on</i>	962
910	<i>shop “Translation Evaluation–From Fragmented</i>	<i>Benchmarking: Past, Present and Future</i> . Associa-	963
911	<i>Tools and Data Sets to an Integrated Ecosystem</i> ,	tion for Computational Linguistics, Online.	964
912	pages 27–34.		
913	Antoine Bosselut, Esin Durmus, Varun Prashant Gan-	Kenneth Ward Church and Joel Hestness. 2019. A sur-	965
914	gal, Sebastian Gehrmann, Yacine Jernite, Laura	vey of 25 years of evaluation. <i>Natural Language</i>	966
915	Perez-Beltrachini, Samira Shaikh, and Wei Xu, edi-	<i>Engineering</i> , 25(6):753–767.	967
916	tors. 2021. <i>Proceedings of the 1st Workshop on Nat-</i>	Sunipa Dev, Emily Sheng, Jieyu Zhao, Jiao Sun,	968
917	<i>ural Language Generation, Evaluation, and Metrics</i>	Yu Hou, Mattie Sanseverino, Jiin Kim, Nanyun	969
918	<i>(GEM 2021)</i> . Association for Computational Lin-	Peng, and Kai-Wei Chang. 2021. What do bias mea-	970
919	guistics, Online.	asures measure? <i>arXiv preprint arXiv:2108.03362</i> .	971
920	Samuel R. Bowman and George Dahl. 2021. What will	Chris Drummond. 2009. Replicability is not repro-	972
921	it take to fix benchmarking in natural language un-	ducibility: nor is it good science. In <i>Proceedings of</i>	973
922	derstanding? In <i>Proceedings of the 2021 Confer-</i>	<i>the Evaluation Methods for Machine Learning Work-</i>	974
923	<i>ence of the North American Chapter of the Associa-</i>	<i>shop at the 26th ICML</i> , volume 1. Citeseer.	975
924	<i>tion for Computational Linguistics: Human Lan-</i>	Jesse Dunietz, Greg Burnham, Akash Bharadwaj,	976
925	<i>guage Technologies</i> , pages 4843–4855, Online. As-	Owen Rambow, Jennifer Chu-Carroll, and Dave Fer-	977
926	sociation for Computational Linguistics.	rucci. 2020. To test machine comprehension, start	978
927	António Branco, Nicoletta Calzolari, and Khalid	by defining comprehension . In <i>Proceedings of the</i>	979
928	Choukri. 2016. Workshop on research results repro-	<i>58th Annual Meeting of the Association for Compu-</i>	980
929	ducibility and resources citation in science and tech-	<i>tational Linguistics</i> , pages 7839–7859, Online. As-	981
930	nology of language. <i>European Language Resources</i>	sociation for Computational Linguistics.	982
931	<i>Association</i> .		
932	Andrea Brennen. 2020. What do people really want	Yoann Dupont. 2019. Un corpus libre, évolutif et ver-	983
933	when they say they want" explainable ai?" we asked	sionné en entités nommées du français (a free, evol-	984
934	60 stakeholders. In <i>Extended Abstracts of the 2020</i>	ing and versioned french named entity recognition	985
935	<i>CHI Conference on Human Factors in Computing</i>	corpus) . In <i>Actes de la Conférence sur le Traitement</i>	986
936	<i>Systems</i> , pages 1–7.	<i>Automatique des Langues Naturelles (TALN) PFIA</i>	987
937	Sabine Buchholz and Erwin Marsi. 2006. CoNLL-	<i>2019. Volume II : Articles courts</i> , pages 437–446,	988
938	X shared task on multilingual dependency parsing .	Toulouse, France. ATALA.	989
939	In <i>Proceedings of the Tenth Conference on Com-</i>	Steffen Eger, Yang Gao, Maxime Peyrard, Wei Zhao,	990
940	<i>putational Natural Language Learning (CoNLL-X)</i> ,	and Eduard Hovy, editors. 2020. <i>Proceedings of</i>	991
941	pages 149–164, New York City. Association for	<i>the First Workshop on Evaluation and Comparison</i>	992
942	Computational Linguistics.	<i>of NLP Systems</i> . Association for Computational Lin-	993
		guistics, Online.	994
		ELSE. 1998. Towards a european evaluation infrastruc-	995
		ture for nl and speech. <i>Workshop at LREC</i> .	996

997	Carla Parra Escartín, Teresa Lynn, Joss Moorkens, and Jane Dunne. 2021. Towards transparency in nlp shared tasks. <i>arXiv preprint arXiv:2105.05020</i> .	1050
998		1051
999		1052
1000	Kawin Ethayarajh and Dan Jurafsky. 2020. Utility is in the eye of the user: A critique of NLP leaderboards . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4846–4853, Online. Association for Computational Linguistics.	1053
1001		1054
1002		1055
1003		1056
1004		1057
1005		1058
1006		1059
1007	Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models . <i>Transactions of the Association for Computational Linguistics</i> , 8:34–48.	1060
1008		1061
1009		1062
1010	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain . In <i>Proceedings of the Sixth Conference on Machine Translation</i> , pages 733–774, Online. Association for Computational Linguistics.	1063
1011		1064
1012		1065
1013		1066
1014		1067
1015		1068
1016		1069
1017		1070
1018		1071
1019	Bin Fu, Yunqi Qiu, Chengguang Tang, Yang Li, Haiyang Yu, and Jian Sun. 2020. A survey on complex question answering over knowledge base: Recent advances and challenges. <i>arXiv preprint arXiv:2007.13069</i> .	1072
1020		1073
1021		1074
1022		1075
1023	Wee Chung Gan and Hwee Tou Ng. 2019. Improving the robustness of question answering systems to question paraphrasing . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6065–6075, Florence, Italy. Association for Computational Linguistics.	1076
1024		1077
1025		1078
1026		1079
1027		1080
1028		1081
1029	Matt Gardner, Jonathan Berant, Hannaneh Hajishirzi, Alon Talmor, and Sewon Min. 2019. Question answering is a format; when is it useful? <i>arXiv preprint arXiv:1909.11291</i> .	1082
1030		1083
1031		1084
1032		1085
1033	Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. AllenNLP: A deep semantic natural language processing platform . In <i>Proceedings of Workshop for NLP Open Source Software (NLP-OSS)</i> , pages 1–6, Melbourne, Australia. Association for Computational Linguistics.	1086
1034		1087
1035		1088
1036		1089
1037		1090
1038		1091
1039		1092
1040		1093
1041	Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. <i>Communications of the ACM</i> , 64(12):86–92.	1094
1042		1095
1043		1096
1044		1097
1045		1098
1046	Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. 2022. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text. <i>arXiv preprint arXiv:2202.06935</i> .	1099
1047		1100
1048		1101
1049		1102
		1103
		1104
	Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. 2018. Explaining explanations: An overview of interpretability of machine learning. In <i>2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)</i> , pages 80–89. IEEE.	1050
		1051
		1052
		1053
		1054
		1055
	Ralph Grishman and Beth Sundheim. 1996. Design of the muc-6 evaluation. Technical report, NEW YORK UNIV NY DEPT OF COMPUTER SCIENCE.	1056
		1057
		1058
		1059
	Mika Härmäläinen and Khalid Alnajjar. 2021. The great misalignment problem in human evaluation of NLP methods . In <i>Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)</i> , pages 69–74, Online. Association for Computational Linguistics.	1060
		1061
		1062
		1063
		1064
		1065
	High-Level Expert Group on AI. 2019. Ethics guidelines for trustworthy ai .	1066
		1067
	Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing.	1068
		1069
		1070
		1071
	Dirk Hovy and Shannon L. Spruit. 2016. The social impact of natural language processing . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 591–598, Berlin, Germany. Association for Computational Linguistics.	1072
		1073
		1074
		1075
		1076
		1077
	Eduard Hovy and Julia Lavid. 2010. Towards a ‘science’ of corpus annotation: a new methodological challenge for corpus linguistics. <i>International journal of translation</i> , 22(1):13–36.	1078
		1079
		1080
		1081
	David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. 2020. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions . In <i>Proceedings of the 13th International Conference on Natural Language Generation</i> , pages 169–182, Dublin, Ireland. Association for Computational Linguistics.	1082
		1083
		1084
		1085
		1086
		1087
		1088
		1089
		1090
		1091
	Yu-Lun Hsieh, Minhao Cheng, Da-Cheng Juan, Wei Wei, Wen-Lian Hsu, and Cho-Jui Hsieh. 2019. On the robustness of self-attentive models . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1520–1529, Florence, Italy. Association for Computational Linguistics.	1092
		1093
		1094
		1095
		1096
		1097
		1098
	Nancy Ide, Nicoletta Calzolari, Judith Ecker-Köhler, Dafydd Gibbon, Sebastian Hellmann, Kiyong Lee, Joakim Nivre, and Laurent Romary. 2017. Community standards for linguistically-annotated resources. In <i>Handbook of linguistic annotation</i> , pages 113–165. Springer.	1099
		1100
		1101
		1102
		1103
		1104

1105	Nancy Ide, Keith Suderman, James Pustejovsky, Marc	Constantine Lignos and Marjan Kamyab. 2020. If you	1161
1106	Verhagen, and Christopher Cieri. 2016. The lan-	build your own NER scorer, non-replicable results	1162
1107	guage application grid and galaxy . In <i>Proceedings</i>	will come . In <i>Proceedings of the First Workshop on</i>	1163
1108	<i>of the Tenth International Conference on Language</i>	<i>Insights from Negative Results in NLP</i> , pages 94–99,	1164
1109	<i>Resources and Evaluation (LREC’16)</i> , pages 457–	Online. Association for Computational Linguistics.	1165
1110	462, Portorož, Slovenia. European Language Re-		
1111	sources Association (ELRA).		
1112	Nancy Ide, Keith Suderman, Marc Verhagen, and	Anne-Laure Ligozat and Sasha Luccioni. 2021. A prac-	1166
1113	James Pustejovsky. 2015. The language application	tical guide to quantifying carbon emissions for ma-	1167
1114	grid web service exchange vocabulary. In <i>Internat-</i>	chine learning researchers and practitioners. Techni-	1168
1115	<i>ional Workshop on Worldwide Language Service In-</i>	cal report, MILA; LISN.	1169
1116	<i>frastructure</i> , pages 18–32. Springer.		
1117	Richard Johansson and Pierre Nugues. 2006. Invest-	Ana Lucic, Maurits Bleeker, Samarth Bhargav, Jessica	1170
1118	tigating multilingual dependency parsing . In <i>Pro-</i>	Forde, Koustuv Sinha, Jesse Dodge, Sasha Luccioni,	1171
1119	<i>ceedings of the Tenth Conference on Computational</i>	and Robert Stojnic. 2022. Towards reproducible ma-	1172
1120	<i>Natural Language Learning (CoNLL-X)</i> , pages 206–	chine learning research in natural language process-	1173
1121	210, New York City. Association for Computational	ing . In <i>Proceedings of the 60th Annual Meeting of</i>	1174
1122	Linguistics.	<i>the Association for Computational Linguistics: Tu-</i>	1175
		<i>torial Abstracts</i> , pages 7–11, Dublin, Ireland. Asso-	1176
		ciation for Computational Linguistics.	1177
1123	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika	Muhammad Kamran Malik and Syed Mansoor Sarwar.	1178
1124	Bali, and Monojit Choudhury. 2020. The state and	2016. Named entity recognition system for postposi-	1179
1125	fate of linguistic diversity and inclusion in the NLP	itional languages: urdu as a case study. <i>International</i>	1180
1126	world . In <i>Proceedings of the 58th Annual Meet-</i>	<i>Journal of Advanced Computer Science and Appli-</i>	1181
1127	<i>ing of the Association for Computational Linguistics</i> ,	<i>cations</i> , 7(10).	1182
1128	pages 6282–6293, Online. Association for Computa-		
1129	tional Linguistics.	Christopher Manning, Mihai Surdeanu, John Bauer,	1183
		Jenny Finkel, Steven Bethard, and David McClosky.	1184
1130	Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh	2014. The Stanford CoreNLP natural language pro-	1185
1131	Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vid-	cessing toolkit . In <i>Proceedings of 52nd Annual</i>	1186
1132	gen, Grusha Prasad, Amanpreet Singh, Pratik Ring-	<i>Meeting of the Association for Computational Lin-</i>	1187
1133	shia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel,	<i>guistics: System Demonstrations</i> , pages 55–60, Bal-	1188
1134	Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mo-	timore, Maryland. Association for Computational	1189
1135	hit Bansal, Christopher Potts, and Adina Williams.	Linguistics.	1190
1136	2021. Dynabench: Rethinking benchmarking in	Benjamin Marie, Atsushi Fujita, and Raphael Rubino.	1191
1137	NLP . In <i>Proceedings of the 2021 Conference of</i>	2021. Scientific credibility of machine translation	1192
1138	<i>the North American Chapter of the Association for</i>	research: A meta-evaluation of 769 papers . In <i>Pro-</i>	1193
1139	<i>Computational Linguistics: Human Language Tech-</i>	<i>ceedings of the 59th Annual Meeting of the Associa-</i>	1194
1140	<i>nologies</i> , pages 4110–4124, Online. Association for	<i>tion for Computational Linguistics and the 11th In-</i>	1195
1141	Computational Linguistics.	<i>ternational Joint Conference on Natural Language</i>	1196
		<i>Processing (Volume 1: Long Papers)</i> , pages 7297–	1197
1142	Jin-Dong Kim, Nancy Ide, and Keith Suderman. 2020.	7306, Online. Association for Computational Lin-	1198
1143	Towards standardization of web service protocols	guistics.	1199
1144	for NLPaaS . In <i>Proceedings of the 1st Inter-</i>	Mónica Marrero, Julián Urbano, Sonia Sánchez-	1200
1145	<i>national Workshop on Language Technology Plat-</i>	Cuadrado, Jorge Morato, and Juan Miguel Gómez-	1201
1146	<i>forms</i> , pages 59–65, Marseille, France. European	Berbis. 2013. Named entity recognition: fallacies,	1202
1147	Language Resources Association.	challenges and opportunities. <i>Computer Standards</i>	1203
		<i>& Interfaces</i> , 35(5):482–489.	1204
1148	Mike Kroutikov. 2019. 7776 ways to compute f1 for an	Yann Mathet, Antoine Widlöcher, Karën Fort, Claire	1205
1149	ner task .	François, Olivier Galibert, Cyril Grouin, Juliette	1206
		Kahn, Sophie Rosset, and Pierre Zweigenbaum.	1207
1150	Jochen L. Leidner and Vassilis Plachouras. 2017. Eth-	2012. Manual corpus annotation: Giving meaning	1208
1151	ical by design: Ethics best practices for natural lan-	to the evaluation metrics . In <i>Proceedings of COL-</i>	1209
1152	guage processing . In <i>Proceedings of the First ACL</i>	<i>ING 2012: Posters</i> , pages 809–818, Mumbai, India.	1210
1153	<i>Workshop on Ethics in Natural Language Process-</i>	The COLING 2012 Organizing Committee.	1211
1154	<i>ing</i> , pages 30–40, Valencia, Spain. Association for		
1155	Computational Linguistics.	Nitika Mathur, Timothy Baldwin, and Trevor Cohn.	1212
1156	Brian Lester. 2020. iobes: Library for span level pro-	2020. Tangled up in BLEU: Reevaluating the eval-	1213
1157	cessing . In <i>Proceedings of Second Workshop for</i>	uation of automatic machine translation evaluation	1214
1158	<i>NLP Open Source Software (NLP-OSS)</i> , pages 115–	metrics . In <i>Proceedings of the 58th Annual Meet-</i>	1215
1159	119, Online. Association for Computational Linguis-	<i>ing of the Association for Computational Linguistics</i> ,	1216
1160	tics.	pages 4984–4997, Online. Association for Computa-	1217
		tional Linguistics.	1218

1219	Kevin McTait and Khalid Choukri. 2003. Setting up an evaluation infrastructure for human language technologies in Europe . In <i>Proceedings of the EACL 2003 Workshop on Evaluation Initiatives in Natural Language Processing: are evaluation methods, metrics and resources reusable?</i> , page 7377, Columbus, Ohio. Association for Computational Linguistics.	1277
1220		1278
1221		
1222		
1223		
1224		
1225		
1226	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In <i>Proceedings of the conference on fairness, accountability, and transparency</i> , pages 220–229.	
1227		
1228		
1229		
1230		
1231		
1232	Mathieu Morey, Philippe Muller, and Nicholas Asher. 2017. How much progress have we made on RST discourse parsing? a replication study of recent results on the RST-DT . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1319–1324, Copenhagen, Denmark. Association for Computational Linguistics.	
1233		
1234		
1235		
1236		
1237		
1238		
1239		
1240	Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. Universal Dependencies v1: A multilingual treebank collection . In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)</i> , pages 1659–1666, Portorož, Slovenia. European Language Resources Association (ELRA).	
1241		
1242		
1243		
1244		
1245		
1246		
1247		
1248		
1249		
1250	Pablo Gamallo Otero and Isaac González. 2012. Dep-pattern: a multilingual dependency parser. In <i>Proceedings of PROPOR</i> . Citeseer.	
1251		
1252		
1253	Chester Palen-Michel, Nolan Holley, and Constantine Lignos. 2021. SeqScore: Addressing barriers to reproducible named entity recognition evaluation . In <i>Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems</i> , pages 40–50, Punta Cana, Dominican Republic. Association for Computational Linguistics.	
1254		
1255		
1256		
1257		
1258		
1259		
1260	Pierre-Henri Paris and Fabian Suchanek. 2021. Non-named entities—the silent majority. In <i>European Semantic Web Conference</i> , pages 131–135. Springer.	
1261		
1262		
1263	Rebecca J. Passonneau and Bob Carpenter. 2014. The benefits of a model of annotation . <i>Transactions of the Association for Computational Linguistics</i> , 2:311–326.	
1264		
1265		
1266		
1267	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.	
1268		
1269		
1270		
1271		
1272		
1273	Barbara Plank, Dirk Hovy, and Anders Søgaard. 2014. Linguistically debatable or just plain wrong? In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 507–511, Baltimore, Maryland. Association for Computational Linguistics.	1277
1274		1278
1275		
1276		
	Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation . In <i>Proceedings of the Tenth Workshop on Statistical Machine Translation</i> , pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.	1279
		1280
		1281
		1282
		1283
	Matt Post. 2018. A call for clarity in reporting BLEU scores . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	1284
		1285
		1286
		1287
		1288
	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	1289
		1290
		1291
		1292
		1293
		1294
	Georg Rehm, Maria Berger, Ela Elsholz, Stefanie Hegele, Florian Kintzel, Katrin Marheinecke, Stelios Piperidis, Miltos Deligiannis, Dimitris Galanis, Katerina Gkirtzou, Penny Labropoulou, Kalina Bontcheva, David Jones, Ian Roberts, Jan Hajič, Jana Hamrlová, Lukáš Kačena, Khalid Choukri, Victoria Arranz, Andrejs Vasiljevs, Oriens Anvari, Andis Lagzdīņš, Jūlija Melņika, Gerhard Backfried, Erīņ Dikīci, Miroslav Janosik, Katja Prinz, Christoph Prinz, Severin Stampler, Dorothea Thomas-Aniola, José Manuel Gómez-Pérez, Andres Garcia Silva, Christian Berrío, Ulrich Germann, Steve Renals, and Ondrej Klejch. 2020a. European language grid: An overview . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 3366–3380, Marseille, France. European Language Resources Association.	1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
	Georg Rehm, Dimitris Galanis, Penny Labropoulou, Stelios Piperidis, Martin Weiß, Ricardo Usbeck, Joachim Köhler, Miltos Deligiannis, Katerina Gkirtzou, Johannes Fischer, Christian Chiarcos, Nils Feldhus, Julian Moreno-Schneider, Florian Kintzel, Elena Montiel, Víctor Rodríguez Doncel, John Philip McCrae, David Laqua, Irina Patricia Theile, Christian Dittmar, Kalina Bontcheva, Ian Roberts, Andrejs Vasiljevs, and Andis Lagzdīņš. 2020b. Towards an interoperable ecosystem of AI and LT platforms: A roadmap for the implementation of different levels of interoperability . In <i>Proceedings of the 1st International Workshop on Language Technology Platforms</i> , pages 96–107, Marseille, France. European Language Resources Association.	1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
	Georg Rehm, Stelios Piperidis, Kalina Bontcheva, Jan Hajic, Victoria Arranz, Andrejs Vasiljevs, Gerhard Backfried, Jose Manuel Gomez-Perez, Ulrich Germann, Rémi Calizzano, Nils Feldhus, Stefanie Hegele, Florian Kintzel, Katrin Marheinecke, Julian Moreno-Schneider, Dimitris Galanis, Penny Labropoulou, Miltos Deligiannis, Katerina Gkirtzou,	1328
		1329
		1330
		1331
		1332
		1333
		1334

1335	Athanasia Kolovou, Dimitris Gkoumas, Leon Voukoutis, Ian Roberts, Jana Hamrlova, Dusan Varis, Lukas Kacena, Khalid Choukri, Valérie Mapelli, Mickaël Rigault, Julija Melnika, Miro Janosik, Katja Prinz, Andres Garcia-Silva, Cristian Berrio, Ondrej Klejch, and Steve Renals. 2021. European language grid: A joint platform for the European language technology community . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations</i> , pages 221–230, Online. Association for Computational Linguistics.	1393
1336		1394
1337		1395
1338		
1339		1396
1340		1397
1341		1398
1342		1399
1343		
1344		1400
1345		1401
1346		1402
1347	Ehud Reiter. 2018. A structured review of the validity of BLEU . <i>Computational Linguistics</i> , 44(3):393–401.	1403
1348		1404
1349		
1350		1405
1351	Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of NLP models with CheckList . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4902–4912, Online. Association for Computational Linguistics.	1406
1352		1407
1353		1408
1354		1409
1355		1410
1356		1411
1357		
1358	Anna Rogers. 2019. How the transformers broke nlp leaderboards. <i>Posted on the Hacking Semantics blog</i> : https://hackingsemantics.xyz/2019/leaderboards .	1412
1359		1413
1360		1414
1361		1415
1362	Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2021. Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension. <i>arXiv preprint arXiv:2107.12708</i> .	1416
1363		1417
1364		
1365	Marta Sabou, Kalina Bontcheva, Leon Derczynski, and Arno Scharl. 2014. Corpus annotation through crowdsourcing: Towards best practice guidelines . In <i>Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)</i> , pages 859–866, Reykjavik, Iceland. European Language Resources Association (ELRA).	1418
1366		1419
1367		1420
1368		1421
1369		1422
1370		1423
1371		1424
1372	Ivan Sanchez, Jeff Mitchell, and Sebastian Riedel. 2018. Behavior analysis of NLI models: Uncovering the influence of three factors on robustness . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1975–1985, New Orleans, Louisiana. Association for Computational Linguistics.	1425
1373		1426
1374		1427
1375		1428
1376		
1377		1429
1378		1430
1379		1431
1380		1432
1381		1433
1382	Michael Saxon, Sharon Levy, Xinyi Wang, Alon Albalak, and William Yang Wang. 2021. Modeling disclosive transparency in NLP application descriptions . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2023–2037, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1434
1383		1435
1384		1436
1385		1437
1386		1438
1387		1439
1388		1440
1389		1441
1389	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.	1442
1390		1443
1391		1444
1392		1445
		1446
	Wei Shen, Jianyong Wang, and Jiawei Han. 2014. Entity linking with a knowledge base: Issues, techniques, and solutions. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 27(2):443–460.	
	György Szarvas, Richárd Farkas, and András Kocsor. 2006. A multilingual named entity recognition system using boosting and c4.5 decision tree learning algorithms. In <i>International Conference on Discovery Science</i> , pages 267–278. Springer.	
	Bruno Taillé, Vincent Guigue, Geoffrey Scoutheeten, and Patrick Gallinari. 2020. Let's Stop Incorrect Comparisons in End-to-end Relation Extraction! In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3689–3701, Online. Association for Computational Linguistics.	
	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition . In <i>Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003</i> , pages 142–147.	
	Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander Wubben, and Emiel Kraemer. 2019. Best practices for the human evaluation of automatically generated text . In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pages 355–368, Tokyo, Japan. Association for Computational Linguistics.	
	Christopher Walker, Stephanie Strassel, Julie Medero, and Kazuaki Maeda. 2006. Ace 2005 multilingual training corpus. <i>Linguistic Data Consortium, Philadelphia</i> , 57:45.	
	Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. CharacTer: Translation edit rate on character level . In <i>Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers</i> , pages 505–510, Berlin, Germany. Association for Computational Linguistics.	
	Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al. 2016. The fair guiding principles for scientific data management and stewardship. <i>Scientific data</i> , 3(1):1–9.	
	Billy T.M. Wong and Sophia Y.M. Lee. 2013. Annotating legitimate disagreement in corpus construction . In <i>Proceedings of the 11th Workshop on Asian Language Resources</i> , pages 51–57, Nagoya, Japan. Asian Federation of Natural Language Processing.	

- 1447 Peng Xu and Frederick Jelinek. 2004. [Random forests](#)
1448 [in language modelin](#). In *Proceedings of the 2004*
1449 *Conference on Empirical Methods in Natural Lan-*
1450 *guage Processing*, pages 325–332, Barcelona, Spain.
1451 Association for Computational Linguistics.
- 1452 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q
1453 Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-
1454 uating text generation with bert. *arXiv preprint*
1455 *arXiv:1904.09675*.
- 1456 Shuyan Zhou, Xiangkai Zeng, Yingqi Zhou, Antonios
1457 Anastasopoulos, and Graham Neubig. 2019. [Im-](#)
1458 [proving robustness of neural machine translation](#)
1459 [with multi-task learning](#). In *Proceedings of the*
1460 *Fourth Conference on Machine Translation (Volume*
1461 *2: Shared Task Papers, Day 1)*, pages 565–571, Flo-
1462 rence, Italy. Association for Computational Linguis-
1463 tics.
- 1464 Thomas Zielke. 2020. Is artificial intelligence ready
1465 for standardization? In *European Conference*
1466 *on Software Process Improvement*, pages 259–274.
1467 Springer.