

LLMs in the Real World: Evaluating “AI” in Emergency Contexts

Anonymous ACL submission

Abstract

We present a case study on the initial stages of an LLM-based machine translation application’s deployment in a real-world context: a text-2-911 system advertising capabilities in 55 languages for use in emergencies in which it may be difficult to call operators directly. We identify a number of common misconceptions about technologies such as these and describe their implications, concluding with a set of concrete recommendations and best practices for stakeholders at every stage of the development and deployment pipeline. We offer a call to action and urge our colleagues in the research community to play a greater role in the articulation of our findings to the public. While the advancement of scientific research often lies in solving the “hard” problems, we argue that it is often the “easy” ones— problems for which the latest technology is often unnecessary— that are most overlooked.

1 Introduction

Despite considerable overlap between academic and industry-based developers of Large Language Models (LLMs) and related technologies (Abdalla et al., 2023), it seems Natural Language Processing (NLP) researchers have a science outreach problem. As our research findings continue to drive the development of some of the most quickly adopted user-facing applications to date (De Bruger, 2023), the findings themselves— and their real-world implications— are too often lost in the hype. Artificial Intelligence (AI) is increasingly being offered as an inevitable solution to some of humanity’s largest problems (Eubanks, 2018; Byrum and Benjamin, 2022; Benjamin, 2024; Center for Democracy & Technology, 2025b).

As the outputs of modern NLP research are taken up and applied in commercial products, an information gap has developed between those designing NLP applications and their end users. For example,

researchers may take it for granted that a model performs worse with lower-resourced languages (Silva et al., 2024), or that its performance degrades in response to inputs from outside its training domain (Wu et al., 2024; Li et al., 2025). However, decision makers deploying NLP products in essential services like law enforcement and emergency response may not be aware of these limitations. Many may find themselves under pressure to find ways to acquire and integrate AI tools, but are left to navigate the AI software market without the knowledge needed to properly evaluate them, mitigate their risks, and ensure that they’re as safe, ethical, and effective as possible.

This failure to effectively communicate our research findings to the public— including to those developing and selling consumer-facing applications— is not unique to NLP. Cryptographers have developed easy-to-access, publicly available code for secure public-key communication, but surveys of these systems in actual use persistently reveal that end users continue to create security flaws by misusing the APIs (Choudhari et al., 2021; Lazar et al., 2014). The NLP research community now faces a similar problem. Although basic tools for transparency and evaluation like model cards (Mitchell et al., 2019) have existed for years, they still aren’t widely used in commercial settings.

Most consumers of NLP technologies have little other than promotional materials to go off of as they navigate the complex landscape of tools marketed to the public as “AI.” The result, we claim, is both an information imbalance and an accountability gap: NLP products are increasingly being sold to consumers in both public and private sectors, including in high stakes contexts such as policing (United States v. Cruz-Zamora, 2018; Quaglia, 2022), immigration court (Deck, 2023), critical public health announcements (Moreno, 2021), and other emergency response (Burns, 2025), without

adequate information or support to use them safely. Furthermore, if harm occurs as a result of technological error or misuse, it is often unclear who, if anyone, can be held accountable for that harm. We believe this should be cause for genuine concern within the research community. When language technology is deployed to support emergency services within our own communities, the stakes are high and all of us are stakeholders.

This paper presents a case study on one example of an LLM-based language technology being bought, sold, and deployed for use in emergencies at a local 911 center. We describe the technology’s rollout via its marketing and promotional materials, and describe our experience meeting with staff involved in its deployment at the 911 center.

Our experience sheds light on a number of common misconceptions about language and language technology that, in combination with systemic gaps in accountability, can result in the risky and potentially harmful deployment of NLP systems. This is of particular concern in situations such as our case study, in which the product is being used in emergencies and affects some of the most vulnerable members of our community— refugees and other immigrants for whom English is not a native language. We discuss the role that NLP researchers can play in addressing these critical gaps, and conclude the paper with a set of concrete recommendations for best practices. Finally, we encourage our colleagues in the research community to do more to support science outreach, so that our advice may be audible to those who need to hear it.

2 Case Study: Text-2-911 Service

2.1 Background

According to U.S. federal and state statutes, including Title IV of the Civil Rights Act of 1964, the Americans with Disabilities Act, the Affordable Care Act, and the 14th Amendment, among others, emergency service providers are legally obligated to ensure language access for callers with limited English proficiency. While the ability to communicate with emergency responders may be assumed as a given if one’s native language is English, lack of equitable language access can compound hardships faced by many of our most vulnerable populations, including immigrant and refugee communities, as well as individuals with disabilities (National Immigrant Women’s Advocacy Project (NiWAP) and American University Washington College of Law,

2013; Taira et al., 2021; Bhuiyan, 2023; Hofmann et al., 2024; Parmar, 2025). Community organizations can provide crucial support— sometimes the only support— for such individuals, helping them to navigate complex systems and institutions with which they may be unfamiliar, educating them about their rights, and connecting them to critical resources.

The second author of this study is a licensed social worker at one such community organization, leading a multilingual team of victim advocates specialized in serving immigrant and refugee survivors with limited English proficiency. Through this work, she and her colleagues have seen firsthand how a lack of qualified interpreters and misuse of language technology can lead to a cascade of harmful effects on vulnerable populations. So when the city announced the rollout of a new AI-powered¹ translator for text-2-911 emergency response, she and her team were eager to learn more.

The option to text 911 in English began as a solution for Deaf and Hard of Hearing callers and “those who may be unable to communicate verbally due to background noise or safety considerations” (Laird, 2025) and became available in early 2019. Crucially, the text-2-911 tool is not intended to be an equivalent alternative to a voice call (Franklin County Board of Commissioners, 2019). In our visits to the 911 center, described in Section 2.3, the staff made clear that voice calls are always preferred since they can provide additional information, such as background noise and voice distress, to the call operator. Dispatchers are trained to ask anyone who texts the 911 center whether they are able to call safely, and encourage them to do so if possible. The city already provides human interpreters for voice calls to serve the estimated 6.4% of area residents over 5 years old who speak English “less than very well” (Central Ohio Hospital Council et al., 2025), and interpreter services were used for about 4,000 of the 670,000 total calls to the local 911 call center in 2024 (Laird, 2025).

2.2 Popular Perceptions of Machine Translation

More often than not, press coverage of MT applications portrays systems in an overwhelmingly positive light, with relatively little scrutiny of the limitations of the technology or its implications

¹Despite repeated attempts to determine the exact model architecture powering the service, we are still not 100% certain of its design.

for human interpretation (Vieira et al., 2021). The deployment of MT systems, even in high stakes contexts, is often contrasted with the option of providing no language access rather than compared to the statutory baseline, namely in-person or teleinterpretation by qualified, human translators (Quaglia, 2022). When the bar is set artificially low, it is easier for well-intentioned community members to celebrate the use of MT as “above and beyond” when it may actually represent a step backward in access for callers with limited English proficiency.

In our case, one local media report suggests that fear of language barriers when calling 911 “may soon be a thing of the past.” The article goes on to state that “instead of relying on language interpreters to help non-English speaking callers [...] callers can now text 911 in their own language,” contradicting the stated intent of the service not to supplant voice-call interpretation but rather add an additional accessibility option (Keller, 2025).

City residents interviewed about the new technology for the promotional video echo the language from the original press release, assuring residents that they could now text 911 “in their native language.” However, a list of the 55 languages supported by the model has not been included in any press coverage, and we were only able to acquire it by contacting the 911 center directly. Not all of the county’s most commonly spoken languages are on the list, and non-Latin scripts can only be sent through AT&T, a disclaimer that could be easily missed in much of the press coverage.

It was clear that more information was needed if the victim services team wished to provide accurate information and guidance to their clients. This is when the second author reached out to a local university’s linguistics department for clarity on current MT technology. She also contacted the 911 center staff, who immediately and graciously invited her team to visit the center for a tour and in-depth discussion of the new features. Due to ongoing updates to the software, they were not able to test out the translation on that day, so a second meeting was arranged, and the first author was invited to join. The following section describes what was learned from these meetings.

2.3 Visiting the 911 Center

The first and second authors visited the 911 center at the end of September, 2025, along with two colleagues from the Victim Services Program who were eager to test out the translation tool in their

respective native languages. Each had prepared a list of phrases from real-world text messages to explore how the model handled language-specific challenges such as dialect variation, text speak, typos, referential ambiguity, idioms, and code switching. Three city staff members from the 911 center and a representative from the software company providing the MT application generously made time to host us, even as they orchestrated the day’s emergency response activities for a city of 900,000 residents.

The 911 center is home to the city’s Public Safety Answering Point (PSAP), where operators receive all incoming local calls and texts to 911 before routing them to the appropriate responders, e.g., fire, EMS, or police. Software and maintenance for the PSAP interface, including the text-2-911 feature, is provided by a third party who advertises their use of Microsoft Azure to provide language detection and automatic translation. According to 911 center staff, Microsoft does not provide access to the underlying model or training data.

The staff managing implementation of the tool had also not been provided any evaluation data or quality assurance services by their software provider. While a policy exists at the state level outlining deliberate and detailed requirements for “planning, implementation, procurement, security, privacy, and governance requirements for the use of Artificial Intelligence (AI)” (State of Ohio, 2023), no equivalent policy has been created for City departments. This appears to leave the 911 center without the necessary subject-matter expertise, training, guidance, or resource allocation to ensure that proper safeguards are in place.

According to the software company’s representative, the goal of the translation tool is to decrease response time for end users with limited English proficiency. However, there is currently no ongoing evaluation to establish the product’s success toward meeting that goal. The MT system also does not integrate any oversight from human translators, either in real-time or after-the-fact for quality assurance. A human dispatcher receives and responds to the translated text, but the text output by the MT model, as with any AI model, is still ultimately AI-generated.

2.4 Testing the Tool

We were given the opportunity to interact with the system ourselves in real time. We probed the model with common linguistic phenomena found in text

283 messages, such as accidental misspellings and di- 335
284 alectal variation. Ultimately, both of our colleagues 336
285 from the Victim Services Program encountered 337
286 challenges texting 911 in their native languages, 338
287 Arabic and Nepali, respectively. 339

288 In the case of Arabic, we learned that Modern 340
289 Standard Arabic (MSA) was the only variety of 341
290 Arabic in which the MT system was supposed to 342
291 be able to interact (Al-Laith and Kebdani, 2025; 343
292 Mishra et al., 2025). However, in addition to the 344
293 limitations posed by Arabic’s non-Latin orthogra- 345
294 phy, those familiar with the sociolinguistic contexts 346
295 in which Arabic is spoken will know that MSA is 347
296 rarely, if ever, the language a speaker will use to 348
297 communicate via text message. Dialectal varia- 349
298 tions in lexical items and spelling presented clear 350
299 challenges for the model’s ability to interact with 351
300 an Arabic speaker via text. Similarly, the MT sys- 352
301 tem is only able to recognize Nepali written in the 353
302 Devanagari script. However, at least among the 354
303 Bhutanese-Nepali community making up the ma- 355
304 jority of Nepali speakers in the area, text messages 356
305 are almost exclusively written using the Latin script. 357
306 Our colleague was not even sure how to use a De- 358
307 vanagari keyboard, and could not find all of the 359
308 symbols he would need to interact with the system 360
309 using the orthography that the model was trained 361
310 on for the language. 362

311 The potential negative impact of such incongru- 363
312 ence between the data the model was trained on 364
313 and that which might be encountered in a realistic 365
314 setting is further amplified by the tool’s lack of an 366
315 informed consent procedure for end users. When 367
316 a person tries to contact 911 via text message in 368
317 a language other than English, the interface dis- 369
318 plays both source text and translation output to the 370
319 dispatcher at the 911 center, who typically is not 371
320 proficient in the target language. In contrast, the 372
321 use of MT is not explicitly disclosed and neither the 373
322 translations nor the name of the language detected 374
323 by the model are presented to the person texting 375
324 911. On the surface, their experience is no differ- 376
325 ent than if they were texting directly to another 377
326 person. They receive no disclosure of AI-generated 378
327 translation or user guidance for optimizing output 379
328 accuracy. 380

329 There is a reasonable concern that wordy dis- 381
330 claimers or excessive instruction could slow down 382
331 the response time or cause confusion in an emer- 383
332 gency situation. These are valid considerations that 384
333 warrant empirical investigation. However, lacking 385
334 that evidence, and given the vast amount of evi-

dence of the errors and risks associated with LLM 335
technologies already documented by the research 336
community (Costa et al., 2015; Berk, 2021; Fre- 337
itag et al., 2021; Mehandru et al., 2023; Court and 338
Elsner, 2024; Freitag et al., 2024; Mickus et al., 339
2024; Urlana et al., 2025), we wish to emphasize 340
that requiring informed consent and transparency 341
for users on both sides of the interaction is not only 342
more ethical, it is also likely to improve the tool’s 343
overall performance and make the service more 344
effective. 345

346 During the meeting at the 911 center, we dis- 347
348 cussed a number of additional ethical considera- 348
349 tions and safety precautions. Although most of our 349
350 concerns have been well-documented in the aca- 350
351 demic literature for years (see, e.g., Kumar et al. 351
2023 for an overview), public officials and local 352
353 decision makers are often unaware of many of the 353
354 ethical best practices NLP researchers might now 354
355 take for granted (Karamolegkou et al., 2025). There 355
356 is still a general lack of understanding of the po- 356
357 tential vulnerabilities and risks inherent to these 357
358 technologies, reflecting what we see as an overall 358
359 gap in access to information and insufficient in- 359
360 volvement or support from experts in our field. Ad- 360
361 dressing these issues at all stages of the model’s life 361
362 cycle likely requires action by city leadership, in- 362
363 cluding legislative policy, resource allocation, and 363
364 partnerships with local community organizations 364
and members of the NLP research community. 365

365 2.5 Learning from Experience 365

366 We hope our case study will encourage our col- 366
367 leagues in the research community to play a greater 367
368 role in advocating for the safe and responsible ap- 368
369 plication of their own findings. However, it should 369
370 be noted that although a sizeable number of ACL 370
371 submissions each year are authored or co-authored 371
372 by researchers in the private sector (Abdalla et al., 372
2023), NLP research culture itself has continued to 373
374 shift away from open-source principals and peer- 374
375 reviewed science towards greater secrecy and a 375
376 “move fast and break things” approach that priori- 376
377 tizes profits and tends to benefit only a small subset 377
378 of the global population (Benjamin, 2019; Blodgett 378
et al., 2020; Junker, 2024). 379

380 Furthermore, unlike in the medical field where 380
381 AI adoption has been carefully monitored and regu- 381
382 lated, many of the groups buying and selling these 382
383 technologies in other domains lack the technical 383
384 expertise, clear guidance, or necessary resources to 384
385 audit and evaluate the deployment at the necessary 385

scale (Vieira et al., 2021). In the multitude of situations in which a researcher isn't present to evaluate an AI product with a critical eye, the success of an LLM application's deployment depends heavily on the pre-existing knowledge and abilities of those acquiring and using it.

The following section describes a number of specific misconceptions about language and language technologies that we've repeatedly observed in circulation among the general public. Perpetuated by the broader societal patterns described in Section 4 and without the support of experts in our field to counteract them, we believe these gaps in knowledge will continue to enable instances of inappropriate, ineffective, and sometimes even harmful deployment of LLM-based language technologies.

3 Common Misconceptions about "AI"

Computer science knowledge or AI literacy among those buying, selling, using, and regulating NLP technologies has consistently lagged behind the speed at which the field has advanced. Journalists who may otherwise provide a source of oversight and information are often also un- or under-informed, which can mask important considerations and mislead the general public (Vieira, 2020). It is worth asking where the following misconceptions come from, and we encourage the research community to do more to publicly debunk them.

3.1 Misconception 1: The Term "AI" is Well-Defined

Since its inception as a field of study, there has been debate about what actually constitutes artificial intelligence (Turing, 1950). "AI" has become a catchall term for a wide variety of large pretrained models, many of which are rapidly becoming a part of our everyday landscape. This can generate an unfortunate—and inaccurate—impression of homogeneity, obfuscating the difference between NLP tasks and objectives, such as those involved in machine translation vs. a dialogue system, or between model architectures, such as the distinctions between LLM-based pipelines and traditional NMT. Conflating these systems into the umbrella term "AI" contributes to the belief that all of this technology is the same, with the same capabilities, costs, and problems.

3.2 Misconception 2: AI has Superhuman Intelligence

It is common for even researchers to anthropomorphize language technologies that sometimes display what can feel like superhuman abilities, like recalling specific facts about more than an encyclopedia's worth of topics (Deshpande et al., 2023; Erscoi et al., 2023). Having already passed the Turing test with flying colors for years, models may now be marketed to consumers as possessing or approaching Artificial General Intelligence (AGI)—a markedly superhuman ability whose actual definition is just as vague and debatable as any other kind of intelligence (Mahowald et al., 2024; Mitchell, 2024). Assuming there is no viable alternative, it is understandable that consumers looking to serve speakers of languages other than English might turn to a seemingly superhuman MT system in an effort to provide something rather than nothing. Even with the best of intentions, however, confusion between supporting a language and supporting it well can have dire consequences (Bhuiyan, 2023; CalMatters, 2025; Center for Democracy & Technology, 2025a; Quaglia, 2022; Deck, 2023).

3.3 Misconception 3: Language is Easy

Potential users of "AI" are not just unaware of the fine points of language technology; they often also hold a variety of misconceptions about language itself (Wagner et al., 2023). These issues can be mutually reinforcing—some have told us they want "translation" rather than "interpretation" because they want to know word-for-word what their interlocutor is saying. Linguists and translation theorists know that this isn't the right approach: the individual words don't always communicate the core meaning, and utterances can be ambiguous or multivalent even for an experienced interpreter (Nielsen et al., 2025). But this folk theory of translation contributes to the misguided belief that machine translation is more objective and therefore more accurate than any human interpreter.

People may also hold misconceptions about language diversity, for example assuming that dialectal variation is merely a matter of accent or that stigmatized dialects are language errors resulting from poor education (Hudley et al., 2024). Given such misunderstandings, it may be easy to believe that a technology advertising support for over 50 languages will be able to serve all of them equally and that dialectal variation will not cause significant

482 problems, contrary to findings from NLP research
483 (Aycock et al., 2025; Hofmann et al., 2024).

484 **3.4 Misconception 4: Quantitative Metrics are** 485 **Reliable and Sufficient**

486 As a largely empirical discipline, NLP relies heav-
487 ily on automatic and quantitative metrics. NLP
488 researchers commonly acknowledge the inadequa-
489 cies of their own metrics (Flamich et al., 2025) and
490 may even take part in shared tasks attempting to
491 improve them (Freitag et al., 2024; Shayegh et al.,
492 2025). Unfortunately, awareness of a metric’s limi-
493 tations too often fails to make it beyond academic
494 circles. In contrast, techniques used to market and
495 sell LLM technologies leverage benchmarks to ad-
496 vertise some of the “superhuman” capabilities dis-
497 cussed in Section 3.2. End users may be unaware
498 that even the best performing model will degrade
499 outside its training domain (Saunders, 2022), and
500 benchmark scores can be gamed (Mansurov et al.,
501 2025). Moreover, simply interpreting the metric
502 numbers can be difficult for novices. Long expe-
503 rience of evaluation gives professionals a general
504 notion of how to mentally map between MT metric
505 scores and translation quality (e.g. Scarton et al.,
506 2019) Without this experience, one might incor-
507 rectly assume that a high score means that mission-
508 critical errors have already been eliminated.

509 **3.5 Misconception 5: Technological** 510 **Solutionism**

511 The problem of over-estimating the abilities of tech-
512 nology while underestimating our own is not a new
513 one. “Solutionism” (Morozov, 2013) is the ten-
514 dency to assume that social problems are amenable
515 to engineering solutions— especially quick, cheap
516 and disruptive ones. The misconception is not that
517 technology cannot help; it often can! But in order
518 to do so, it needs to be embedded within a sup-
519 portive social context (Benjamin, 2024; Sanchez
520 et al., 2025). Many of the problems AI is being
521 sold to fix would likely be more efficiently and ef-
522 fectively resolved with simpler methods (Quaglia,
523 2022). Instead, vendors of so-called “AI solutions”
524 often market their products by communicating, di-
525 rectly or indirectly, that the human element can be
526 dispensed with entirely.

527 However, human oversight is essential when de-
528 ploying technologies as erratic, unpredictable, and
529 potentially even deceptive as LLMs (Ouyang et al.,
530 2022; Roose, 2023; Mickus et al., 2024; Center for
531 Democracy & Technology, 2025b). LLM software

532 providers should therefore be expected to provide
533 training and support for human quality assurance
534 teams, as well as ongoing monitoring in collabora-
535 tion with professional human interpreters to sys-
536 tematically collect and review feedback from the
537 app’s end users. The question “what if something
538 goes wrong” is central to engineering robust sys-
539 tems (Kapur et al., 2014, ch. 1.6, ch. 10). Failure
540 to ask and answer this question can make a system
541 appear relatively cheap to deploy, but this is only
542 because its true costs appear primarily in scenar-
543 ios where it *doesn’t* work. The reported cost of
544 responding to a domestic violence homicide, for
545 example, stands in the millions of dollars (Nessen,
546 2025).

547 Such errors are not inevitable, and their related
548 losses could be minimized by prioritizing the train-
549 ing and employment of professional human inter-
550 preters over high-tech solutions, particularly when
551 it is not possible to guarantee the safety of the tech-
552 nology in question. The results of doing so would
553 not only be “better than nothing,” they would be
554 quantifiably better than a machine translation sys-
555 tem of variable or unknown quality. Investment
556 in humans and the things we do best, such as lan-
557 guage and translation, would thus make for a sound
558 financial (as well as ethical) decision.

559 **4 Why the Problem Persists**

560 We believe our case study represents broader trends
561 in the deployment of LLM-based language tech-
562 nologies around the globe. Why does this happen?

563 **4.1 Information Asymmetries**

564 Without a doubt, one of the biggest contributing fac-
565 tors to the inappropriate and sometimes unethical
566 use of LLM-based technologies is that consumers,
567 even with the best of intentions, lack access to in-
568 formation. There is an AI literacy crisis at nearly
569 every level of the adoption chain.

570 Information asymmetries begin before many pre-
571 trained models are even released to the public. In
572 many cases, it is only possible to infer what the
573 model was trained on by considering its outputs in
574 light of the old computer science adage: “garbage
575 in, garbage out.” The details of a model’s develop-
576 ment is further obscured once it is packaged into
577 software and sold by a third party. Those selling
578 the software do not necessarily understand how
579 their product works in technical terms, and even
580 when developers understand the API they are using,

they have not necessarily been trained in the core technologies behind it.

Analogous to the problem faced by cryptographers described in Section 1, the mass availability of APIs for LLMs creates the illusion that no specialized knowledge is needed to use the product. Similar to using a database server or hash function, it’s easy for an engineer to assume that a large company like Microsoft or Meta has made their product available because it “works,” without a full understanding of what “working AI software” means or the inherent risks and potential errors that come with deciding to use it.

The research community is not without our share of responsibility, either. Although we have made much progress towards an agreed-upon set of standards for conducting ethical research, the economic and cultural environments in which this work takes place does not tend to value or support science outreach and communication to the public. Regulators and legislators also often lack AI literacy, and their perception of NLP research is dominated by the perspectives of a small handful of powerful companies (Kang, 2025). Without the infrastructure and support to quickly and directly communicate our findings to the public, academic researchers effectively surrender our ability to speak on behalf of our own science.

4.2 The Accountability Gap

For NLP software development to be both safe and effective, knowledge has to move from the research community through multiple layers of transmission. At each of these layers, there is an “accountability gap” to cross: developers and sales people won’t learn best practices unless they have a good reason to (Eubanks, 2018; Hohenstein and Jung, 2020). Some areas (like medicine and law) have well-established systems of individual and institutional accountability that can be applied to NLP tools, for example by transparently defining regulations, formalizing community norms, and creating community-internal resources for learning about these technologies (Landers and Behrend, 2023). In policing and emergency response, the situation seems far less structured (Taira et al., 2021; Parmar, 2025). Without impartial evaluation based on the technical details of the model and its specific use context, even well-intentioned actors are left without adequate guidance on how to use the technology safely.

5 Recommendations and Best Practices

We wish to reiterate that we believe the text-2-911 service and our experience at the 911 center are representative of broader patterns in AI software deployment, rather than an edge case.

5.1 Deploying Language Technologies in High-Stakes Situations

In an alternative scenario, staff at the call center in our case study would have had a much clearer idea of what sort of product they were getting. This should begin at the point of sale: if a company advertises translation in multiple languages, they should be transparent about the relative performance of each language pair and forthcoming about the potential risks and limitations of their software. Just as it is no longer acceptable to buy packaged food without a nutrition label or drugs without a pharmacist’s consultation, model cards (Mitchell et al., 2019) should be required for all software applications trained using machine learning methods. Similar to a pharmaceutical, model cards should distinguish between on-label and off-label uses and clearly communicate the known, potential, and hypothetical risks of deployment in the particular contexts for which they are intended.

An explicit analysis of model failures should also be an expected part of deciding whether to acquire a new language technology prior to its deployment. This would allow organizations to formulate a contingency plan for any errors they observe in the process or otherwise believe to be probable. Ideally, this would also allow for a more informed and efficient use of public resources, including spending on qualified human interpreters to provide backup for the most heavily used language pairs and sensitive applications. Organizations can make these decisions more responsibly by formulating clear policies and standards before acquiring or using any specific product. For suggested minimal policy recommendations, we refer the reader to the SAFE-AI Task Force Guidance (2024, 2025), which we adapt and present in Appendix A.

In general, we recommend consumers pay more careful attention to matching their deployment context(s) with appropriate levels of technological maturity and reliability. While it might be tempting to reach for the most powerful models for applications serving our most critical use cases, we need to be setting the bar higher, not lower, in such scenarios. The idea is not that we shouldn’t be using LLMs

681 in emergencies, but rather that these models need
682 to be more closely evaluated and monitored before
683 deployment. Perhaps 911 isn't the best call service
684 to pilot automatic translation.

685 Once a system is deployed, it should be contin-
686 ually monitored and evaluated on the real data it
687 faces and its scores should be publicly available.
688 Community partners can be valuable collaborators
689 for this kind of evaluation, since they are likely to
690 be best informed of the actual needs and nuances
691 specific to the populations being served. Their in-
692 volvement may also help to distinguish the errors
693 that really matter from those with less serious con-
694 sequences, allowing emergency service providers
695 to better allocate their limited resources.

696 User interfaces and interactions with these tech-
697 nologies also need to be more immediately trans-
698 parent. Both parties should be able to see how the
699 system translates their messages and confirm that
700 the correct language has been identified. End users
701 are more likely to have enough English proficiency
702 to identify errors than dispatchers with little or no
703 exposure to the target language, as well as prior
704 experience with MT performance in their native
705 language. Increasing transparency allows for more
706 robust informed consent in real time.

707 While appropriate AI governance requires fund-
708 ing beyond the product itself, such expenses are
709 negligible in comparison to the hidden costs, mon-
710 etary or otherwise, of bypassing human oversight
711 and getting it wrong when it really matters.

712 **5.2 How the ACL Community Can Help**

713 As researchers, we should remind ourselves pe-
714 riodically that the complex problems we may be
715 trying to solve are not necessarily the biggest issues
716 that stakeholders still face. Model cards and help
717 lines may not seem cutting-edge, but many of us
718 could also use a reality check: NLP practitioners
719 sometimes over-estimate what's considered com-
720 mon knowledge or how much most people actually
721 understand about language technologies.

722 We also encourage more active contributions to
723 local AI literacy initiatives. That is not to say that
724 every researcher ought to also do science outreach,
725 but that we can all strive to be better about support-
726 ing those among us who do. This can take the form
727 of financial support, but may also just mean ad-
728 vocating for these colleagues within our networks.
729 For example, the ACL could recognize researchers
730 involved in public outreach and safety, as some
731 other professional organizations do (e.g., the Lin-

732 guistics, Language, and the Public Award presented
733 annually by [Linguistic Society of America](#)).

734 Finally, we ought to consider coming to a public
735 consensus on key terminology in our field. NLP
736 researchers know that AI is not one thing, so it's
737 important to communicate this to the public. We
738 can advocate for the use of more specific terminol-
739 ogy when describing these products— for example
740 “LLM chatbot” or “machine-generated translation”—
741 which may help those outside the research commu-
742 nity better distinguish between the various models
743 and their intended uses.

744 **6 Conclusion**

745 This paper presents a case study on a situation we
746 believe to be representative of wider patterns in
747 language technology deployment with the poten-
748 tial to inflict serious, but preventable, harm. When
749 accurate translation can mean life or death, the
750 technology providing it needs to be deployed as
751 ethically and safely as possible. Regulatory legisla-
752 tion and community education are both important,
753 but without active engagement from the research
754 community these strategies are unlikely to be suffi-
755 cient to address the range of issues we've described.
756 There must also be clear mechanisms in place to
757 hold accountable those developing, selling, and pro-
758 viding these technologies once NLP research has
759 moved beyond the theoretical or academic realm.
760 In our opinion, it is both unethical and imprac-
761 tical to place the burden of responsibility on the
762 consumer when deploying LLMs in such critical
763 situations as the one we've described in this study.

764 Our intention in sharing these experiences is not
765 simply to criticize developers or users of NLP ap-
766 plications. Rather, we wish to facilitate open dis-
767 cussion among members of the NLP research com-
768 munity, and across the entire web of stakeholders
769 deploying the technologies our research supports.
770 We hope this may be able to at least mitigate some
771 of the harms that can result from our findings mak-
772 ing their way into society, improving the quality of
773 these applications and increasing their beneficial
774 impact in order to make our communities safer for
775 *everyone*.

776 **7 Limitations**

777 The scope of the current paper is limited to one
778 case study, but we believe the conclusions and rec-
779 ommendations we draw from it apply more broadly.
780 We base our discussion primarily around the lan-

guage used when advertising the MT tool described, as well as our meetings with call center staff, with whom we hope to continue collaborating in order to improve the quality and safety of the services being offered. Not only are so-called “AI” systems themselves relatively new innovations, the software we evaluate as part of our case study has only recently been deployed live. There have yet to be enough interactions with the service, nor have we been invited by the call center or software provider, to conduct the kinds of statistical analyses typically used to validate empirical research in NLP. As we hope to have made clear, describing the technology as “AI-powered” also limits our ability to know or describe the exact model(s) being used or how the system was designed. Finally, ethical considerations, described in the following section, also place rightful limits on the data and methods we might use to further evaluate the tool described in our study.

8 Ethics Statement

This paper addresses life or death services for a segment of our community that is currently among the most vulnerable and the most targeted. In doing so, it was essential for attention to be brought to the risks resulting from a lack of AI literacy in the community, without needing to spotlight individual experiences. No data contained in this paper was obtained, directly, or indirectly, through the provision of services for victims of crime or funded with dollars intended to support those services.

Care was also taken to avoid ascribing bad faith to any of the local stakeholders or to assign blame for the issues described in Section 4, all of which are far from unique to our case study. Every public official, first responder, and community member we spoke to has been open and committed to the goal of increasing language access. This is precisely why we feel an equal sense of responsibility to seek out the role that we, as individuals and as a field, can contribute to the realization of that goal.

References

Mohamed Abdalla, Jan Philip Wahle, Terry Ruas, Aurélie Névéol, Fanny Ducel, Saif Mohammad, and Karen Fort. 2023. [The elephant in the room: Analyzing the presence of big tech in natural language processing research](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13141–

13160, Toronto, Canada. Association for Computational Linguistics. 830
831

Seth Aycok, David Stap, Di Wu, Christof Monz, and Khalil Sima’an. 2025. [Can LLMs Really Learn to Translate a Low-Resource Language from One Grammar Book?](#) In *The Thirteenth International Conference on Learning Representations*. 832
833
834
835
836

R. Benjamin. 2019. *Race After Technology: Abolitionist Tools for the New Jim Code*. Polity Press. 837
838

Ruha Benjamin. 2024. *Imagination: A Manifesto*. W. W. Norton & Company, New York. 839
840

Richard A Berk. 2021. [Artificial Intelligence, Predictive Policing, and Risk Assessment for Law Enforcement](#). *Annual Review of Criminology*, 4(1):209–237. 841
842
843

Johana Bhuiyan. 2023. [Lost in AI translation: Growing reliance on language apps jeopardizes some asylum applications](#). *The Guardian*. 844
845
846

Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. [Language \(Technology\) is Power: A Critical Survey of “Bias” in NLP](#). (arXiv:2005.14050). ArXiv:2005.14050 [cs]. 847
848
849
850

Anna Burns. 2025. [Surrey police shooting death prompts calls for interpreter access](#). *Maple Ridge News*. 851
852
853

Greta Byrum and Ruha Benjamin. 2022. [Disrupting the Gospel of Tech Solutionism to Build Tech Justice](#). *Stanford Social Innovation Review*. 854
855
856

CalMatters. 2025. [Deaf Mongolian Immigrant Held by ICE in California for 4 Months with No Access to Interpreter](#). 857
858
859

Center for Democracy & Technology. 2025a. [Content Moderation in the Global South: A Comparative Study of Four Low-Resource Languages](#). 860
861
862

Center for Democracy & Technology. 2025b. [Humans in the loop](#). Civic tech report, Center for Democracy & Technology. 863
864
865

Central Ohio Hospital Council, Columbus Public Health, and Franklin County Public Health. 2025. [Franklin County HealthMap2025: Community Health Needs Assessment](#). 866
867
868
869

Amit Choudhari, Sylvain Guilley, and Khaled Karray. 2021. [Cryscanner: Finding cryptographic libraries misuse](#). In *2021 8th NAFOSTED Conference on Information and Computer Science (NICS)*, pages 230–235. 870
871
872
873
874

Ângela Costa, Wang Ling, Tiago Luís, Rui Correia, and Luísa Coheur. 2015. [A linguistically motivated taxonomy for machine translation error analysis](#). *Mach. Transl.*, 29(2):127–161. 875
876
877
878

879	Sara Court and Micha Elsner. 2024. Shortcomings of LLMs for low-resource translation: Retrieval and understanding are both the problem . In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 1332–1354, Miami, Florida, USA. Association for Computational Linguistics.	933
880		934
881		935
882		
883		936
884		937
		938
885	William De Brugger. 2023. ChatGPT sets record for fastest growing user base: Analyst note . Accessed October 4, 2025.	939
886		940
887		941
888	Andrew Deck. 2023. AI Translation Is Jeopardizing Afghan Asylum Claims . <i>Rest of World</i> .	942
889		943
		944
890	Ameet Deshpande, Tanmay Rajpurohit, Karthik Narasimhan, and Ashwin Kalyan. 2023. Anthropomorphization of AI: Opportunities and risks . In <i>Proceedings of the Natural Legal Language Processing Workshop 2023</i> , pages 1–7, Singapore. Association for Computational Linguistics.	945
891		946
892		947
893		
894		948
895		949
896	Lelia Erscoi, Annelies Véronique Kleinherenbrink, and Olivia Guest. 2023. Pygmalion displacement: When humanising AI dehumanises women .	950
897		951
898		952
899	Virginia Eubanks. 2018. <i>Automating inequality: How high-tech tools profile, police, and punish the poor</i> . St. Martin’s Press.	953
900		954
901		955
902	Gergely Flamich, David Vilar, Jan-Thorsten Peter, and Markus Freitag. 2025. You cannot feed two birds with one score: the accuracy-naturalness tradeoff in translation . <i>arXiv preprint arXiv:2503.24013</i> .	956
903		957
904		958
905		959
906	Franklin County Board of Commissioners. 2019. Residents can now text-to-911 in an emergency. Press release. Available at: https://www.franklincountyohio.gov/files/assets/public/v/1/emergency-management/documents/text-911-news-release.pdf (accessed [11/20/2025]).	960
907		961
908		962
909		
910		963
911		964
912		965
913	Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, Errors, and Context: A Large-Scale Study of Human Evaluation for Machine Translation . <i>Transactions of the Association for Computational Linguistics</i> , 9:1460–1474.	966
914		967
915		968
916		969
917		
918		970
919	Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-Kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are LLMs breaking MT metrics? results of the WMT24 metrics shared task . In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 47–81, Miami, Florida, USA. Association for Computational Linguistics.	971
920		972
921		973
922		974
923		975
924		976
925		977
926		978
927		979
928		980
929		981
930		982
931	Valentin Hofmann, Pratyusha Ria Kalluri, Dan Jurafsky, and Sharese King. 2024. AI generates covertly racist decisions about people based on their dialect . <i>Nature</i> , 633(8028):147–154. Epub 2024 Aug 28.	983
932		984
		985
		986
		987
	Jess Hohenstein and Malte Jung. 2020. AI as a moral crumple zone: The effects of AI-mediated communication on attribution and trust . 106:106190.	
	Anne H Charity Hudley, Christine Mallinson, and Mary Bucholtz. 2024. <i>Decolonizing linguistics</i> . Oxford University Press.	
	Marie-Odile Junker. 2024. Data-mining and extraction: the gold rush of AI on Indigenous languages . In <i>Proceedings of the Seventh Workshop on the Use of Computational Methods in the Study of Endangered Languages</i> , pages 52–57, St. Julians, Malta. Association for Computational Linguistics.	
	Cecilia Kang. 2025. Trump Unveils Plan to Overhaul A.I. Regulation . <i>The New York Times</i> . Accessed: 2025-11-16.	
	Kailash C. Kapur, Michael Pecht, and Andrew P. Sage. 2014. <i>Reliability engineering</i> . Wiley.	
	Antonia Karamolegkou, Sandrine Schiller Hansen, Ariadni Christopoulou, Filippos Stamatiou, Anne Lauscher, and Anders Søgaard. 2025. Ethical concern identification in NLP: A corpus of ACL Anthology ethics statements . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 11618–11635, Albuquerque, New Mexico. Association for Computational Linguistics.	
	Aliah Keller. 2025. Columbus police break language barriers in emergencies with new tools . <i>Spectrum News 1</i> . Published 5:02 AM ET.	
	Sachin Kumar, Vidhisha Balachandran, Lucille Njoo, Antonios Anastasopoulos, and Yulia Tsvetkov. 2023. Language generation models can cause harm: So what can we do about it? an actionable survey . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3299–3321.	
	Jordan Laird. 2025. Columbus upgrades 911 system with text translation in 55 languages, ‘one-way face-time’ . <i>The Columbus Dispatch</i> .	
	Richard N Landers and Tara S Behrend. 2023. Auditing the AI auditors: A framework for evaluating fairness and bias in high stakes AI predictive models . <i>American Psychologist</i> , 78(1):36.	
	David Lazar, Haogang Chen, Xi Wang, and Nickolai Zeldovich. 2014. Why does cryptographic software fail? a case study and open problems . In <i>Proceedings of 5th Asia-Pacific Workshop on Systems, APSys ’14</i> , New York, NY, USA. Association for Computing Machinery.	
	Bryan Li, Jiaming Luo, Eleftheria Briakou, and Colin Cherry. 2025. Leveraging domain knowledge at inference time for LLM translation: Retrieval versus generation . In <i>Proceedings of the 4th International Workshop on Knowledge-Augmented Methods for Natural</i>	

988	<i>Language Processing</i> , pages 91–106, Albuquerque, New Mexico, USA. Association for Computational Linguistics.	Joseph C. Von Nessen. 2025. <i>The Economic Impact of Intimate Partner Violence in Ohio</i> . Report commissioned by Ohio Domestic Violence Network, released Feb. 24, 2025.	1043 1044 1045 1046
991	Linguistic Society of America. <i>Linguistics, Language, and the Public Award</i> .	Elizabeth Nielsen, Isaac Rayburn Caswell, Jiaming Luo, and Colin Cherry. 2025. <i>Alligators all around: Mitigating lexical confusion in low-resource machine translation</i> . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 206–221.	1047 1048 1049 1050 1051 1052 1053 1054
993	Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2024. <i>Dissociating language and thought in large language models</i> . <i>Trends in cognitive sciences</i> , 28(6):517–540.	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. <i>Training language models to follow instructions with human feedback</i> . <i>Advances in neural information processing systems</i> , 35:27730–27744.	1055 1056 1057 1058 1059 1060
998	Jonibek Mansurov, Akhmed Sakip, and Alham Fikri Aji. 2025. <i>Data laundering: Artificially boosting benchmark results through knowledge distillation</i> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8332–8345, Vienna, Austria. Association for Computational Linguistics.	Tekendra Parmar. 2025. <i>Axon’s Draft One Is Designed to Defy Transparency</i> . <i>Mother Jones</i> . Accessed: 2025-10-20.	1061 1062 1063
1000	Nikita Mehandru, Sweta Agrawal, Yimin Xiao, Ge Gao, Elaine Khoong, Marine Carpuat, and Niloufar Salehi. 2023. <i>Physician detection of clinical harm in machine translation: Quality estimation aids in reliance and backtranslation identifies critical errors</i> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 11633–11647, Singapore. Association for Computational Linguistics.	Sofia Quaglia. 2022. <i>Death by machine translation?</i> <i>Slate</i> . Archived at https://perma.cc/6RD2-3TY3 .	1064 1065
1001	Timothee Mickus, Elaine Zosa, Raul Vazquez, Teemu Vahtola, Jörg Tiedemann, Vincent Segonne, Alessandro Raganato, and Marianna Apidianaki. 2024. <i>SemEval-2024 task 6: SHROOM, a shared-task on hallucinations and related observable overgeneration mistakes</i> . In <i>Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)</i> , pages 1979–1993, Mexico City, Mexico. Association for Computational Linguistics.	Kevin Roose. 2023. <i>Bing’s A.I. Chat Reveals Its Feelings: ‘I Want to Be Alive.’</i> <i>The New York Times</i> . Accessed: 2025-10-19.	1066 1067 1068
1002	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. <i>Model cards for model reporting</i> . In <i>Proceedings of the conference on fairness, accountability, and transparency</i> , pages 220–229.	SAFE-AI Task Force. 2024. <i>Interpreting safe AI task force guidance: AI and interpreting services</i> . Technical report, Stakeholders Advocating for Fair and Ethical AI in Interpreting. Version dated July 1, 2024.	1069 1070 1071 1072
1003	Melanie Mitchell. 2024. <i>The metaphors of artificial intelligence</i> . <i>Science</i> , 386(6723):eadt6140.	SAFE AI Task Force and CoSET. 2025. <i>AI Interpreting Solutions Evaluation Toolkit, Part A: Organization, Implementation and Management</i> . Technical report, SAFE AI Task Force and the Coalition for Sign Language Equity in Technology (CoSET).	1073 1074 1075 1076 1077
1004	Sabrina Moreno. 2021. <i>Virginia Uses Google Translate for COVID Vaccine Information. Here’s How That Magnifies Language Barriers, Misinformation</i> . <i>Richmond Times-Dispatch</i> .	Thomas W Sanchez, Marc Brenman, and Xinyue Ye. 2025. <i>The ethical concerns of artificial intelligence in urban planning</i> . <i>Journal of the American Planning Association</i> , 91(2):294–307.	1078 1079 1080 1081
1005	Evgeny Morozov. 2013. <i>To save everything, click here: The folly of technological solutionism</i> . Public Affairs.	Danielle Saunders. 2022. <i>Domain adaptation and multi-domain adaptation for neural machine translation: A survey</i> . <i>Journal of Artificial Intelligence Research</i> , 75:351–424.	1082 1083 1084 1085
1006	National Immigrant Women’s Advocacy Project (NIWAP) and American University Washington College of Law. 2013. <i>Immigrant and limited english proficient victims’ access to the criminal justice system: The importance of collaboration</i> . Technical report, American University, Washington College of Law.	Scarton Scarton, Mikel L. Forcada, Miquel Esplà-Gomis, and Lucia Specia. 2019. <i>Estimating post-editing effort: a study on human judgements, task-based and reference-based metrics of MT quality</i> . In <i>Proceedings of the 16th International Conference on Spoken Language Translation</i> , Hong Kong. Association for Computational Linguistics.	1086 1087 1088 1089 1090 1091 1092
1007		Behzad Shayegh, Jan-Thorsten Peter, David Vilar, Tobias Domhan, Juraj Juraska, Markus Freitag, and Lili Mou. 2025. <i>Feeding two birds or favoring one? adequacy–fluency tradeoffs in evaluation and meta-evaluation of machine translation</i> . In <i>Proceedings of</i>	1093 1094 1095 1096 1097

1098 *the Tenth Conference on Machine Translation (WMT),*
1099 *Volume 1: Research Papers*, pages 269–285, Miami,
1100 Florida, USA. Association for Computational Lin-
1101 guistics.

1102 Ana Silva, Nikit Srivastava, Tatiana Moteu Ngoli,
1103 Michael Röder, Diego Moussallem, and Axel-
1104 Cyrille Ngonga Ngomo. 2024. Benchmarking low-
1105 resource machine translation systems. In *Proceed-*
1106 *ings of the Seventh Workshop on Technologies for*
1107 *Machine Translation of Low-Resource Languages*
1108 *(LoResMT 2024)*, pages 175–185.

1109 State of Ohio. 2023. [Use of Artificial Intelligence in](#)
1110 [State of Ohio Solutions](#). Administrative policy it-17,
1111 Ohio Department of Administrative Services. Issued
1112 by Kathleen C. Madden, Director.

1113 Breena R. Taira, Valerie Kreger, Amanda Orue, and
1114 Lisa C. Diamond. 2021. [A pragmatic assessment](#)
1115 [of google translate for emergency department in-](#)
1116 [structions](#). *Journal of General Internal Medicine*,
1117 36(11):3361–3365.

1118 Alan M. Turing. 1950. Computing machinery and intel-
1119 ligence. *Mind*, 59(236):433.

1120 United States v. Cruz-Zamora. 2018. United states
1121 vs. omar cruz-zamora. The United States Dis-
1122 trict Court for the District of Kansas. Retrieved
1123 from [https://ecf.ksd.uscourts.gov/cgi-bin/](https://ecf.ksd.uscourts.gov/cgi-bin/show_public_doc?2017cr40100-24)
1124 [show_public_doc?2017cr40100-24](https://ecf.ksd.uscourts.gov/cgi-bin/show_public_doc?2017cr40100-24).

1125 Ashok Urlana, Charaka Vinayak Kumar, Bala Mallikar-
1126 junarao Garlapati, Ajeet Kumar Singh, and Rahul
1127 Mishra. 2025. No size fits all: The perils and pit-
1128 falls of leveraging LLMs vary with company size.
1129 In *Proceedings of the 31st International Conference*
1130 *on Computational Linguistics: Industry Track*, pages
1131 187–203.

1132 Lucas Nunes Vieira. 2020. [Machine translation in the](#)
1133 [news: A framing analysis of the written press](#). *Trans-*
1134 *lation Spaces*, 9(1):98–122.

1135 Lucas Nunes Vieira, Minako O’Hagan, and Carol
1136 O’Sullivan. 2021. [Understanding the societal im-](#)
1137 [pacts of machine translation: A critical review of the](#)
1138 [literature on medical and legal use cases](#). *Informa-*
1139 *tion, Communication & Society*, 24(11):1515–1532.

1140 Laura Wagner, Sumurye Awani, Nikole D Patson, and
1141 Rebekah Stanhope. 2023. To what extent does the
1142 general public endorse language myths? *Language*
1143 *and Linguistics Compass*, 17(3):e12486.

1144 Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can
1145 Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-
1146 Wei Kuo, Nan Guan, and 1 others. 2024. Retrieval-
1147 augmented generation for natural language process-
1148 ing: A survey. *arXiv preprint arXiv:2407.13193*.

A Suggested Policies for Ethical AI Use

Suggested Policies for Ethical Use of AI for Interpreting (Minimum Requirements)

- ✓ **Informed Consent** to accept or decline the use of an AI product
 - ✓ **Opt In/Opt Out** of data collection and storage without penalty
 - ✓ **Ability to Shift** between AI and human interpreting at any point in a timely manner
 - ✓ **User-Friendly Grievance Process** to report errors or harm
 - ✓ **Clear Explanations** of policies regarding privacy, confidentiality, and degrees of AI involvement
 - ✓ **End-User Autonomy** throughout the interpreting process
 - ✓ **Evidence of Improvements** to end-user wellbeing and safety are made available to the public
 - ✓ **Public Transparency** of quality metrics and their results over time
 - ✓ **Accountability and Oversight** to hold providers responsible for any errors or harm
-

Table 1: Ethical principles for the use of AI for interpreting, adapted from [SAFE-AI Task Force \(2024\)](#).