

Predicting Compact Phrasal Rewrites with Large Language Models for Automatic Speech Recognition Post Editing

Anonymous ACL submission

Abstract

Large Language Models (LLMs) excel at rewriting tasks such as text style transfer and grammatical error correction. Although the output in these tasks often significantly overlaps with the input, the decoding cost still increases with output length, regardless of the number of overlaps. By leveraging the overlap between the input and the output, Kaneko and Okazaki (2023) proposed model-agnostic edit span representations to compress the rewrites to save computation. They reported an output length reduction rate of nearly 80% with minimal accuracy impact in four rewriting tasks. In this paper, we propose alternative edit *phrase* representations inspired by phrase-based statistical machine translation. We systematically compare our phrasal representations with their span representation. We apply the LLM rewriting model to the task of Automatic Speech Recognition (ASR) post editing and show that our target-phrase-only edit representation has the best efficiency-accuracy trade-off. On the LibriSpeech test set, our method closes 50-60% of the WER gap between the edit span model and the full rewrite model while losing only 10-20% of the length reduction rate of the edit span model.

1 Introduction

Large Language Models pretrained on vast amount of texts and then fine-tuned, instruction-tuned, or prompted for generation tasks have achieved great success in the past few years (Raffel et al., 2020; Brown et al., 2020; Chowdhery et al., 2022; Anil et al., 2023; OpenAI et al., 2024). These models excel at text rewriting tasks, and text style transfer (Reif et al., 2022) and grammatical error correction (Rothe et al., 2021; Fang et al., 2023) in particular.

But the superior quality of these models comes along with a steep increase in the cost of computation. To enable broad deployment for a large user base, it is crucial to reduce the computational cost while maintaining the accuracy.

One of the common characteristics of the above-mentioned rewriting tasks is that their output often repeats spans of text in the input. Exploiting the common sub-strings between the input and the output can result in more compact representations for the output. LLMs can be fine-tuned on examples that map input to their compact rewrite representations instead of plain rewrites. At inference time, decoding output needs to be composed with the input to expand into complete rewrites. Kaneko and Okazaki (2023) gave one such representation, which is a numerical span indexing into the input sequence followed by a target phrase that will substitute the source phrase in the given span. We propose two new alternative representations. The first one uses a source-target phrase pair to represent each rewrite pattern, analogous to phrase pairs used by phrase-based statistical machine translation (Koehn et al., 2003). The second one only uses a target phrase along with left and right context words that appear in the input. We call the new representations *phrase* representations to distinguish them from the *span* representation of Kaneko and Okazaki (2023).

The clear advantage of compact representations over complete rewrite is that the number of decoding steps, and hence the computational cost of inference, is reduced. Kaneko and Okazaki (2023) reported an output length reduction rate of 80%. The disadvantage of such representations is that decoding errors can also lead to inconsistency with the input sequence in the expansion stage, causing an error propagation effect. For example, using the numerical span representation, if the left index or the right index is off by one, the source phrase to be substituted will also be off by one. The concatenation of the context and the substitution can therefore become disfluent. With phrase representations, context words are provided before and after substitution phrases, which can alleviate the problem of disfluency upon substitution. However, phrase

representations have their own problems too. The predicted context phrase may not match the input, which makes it necessary to discard the subsequent rewrite. The main focus of the paper is to evaluate the efficiency-accuracy trade-off of the different representations.

We choose Automatic Speech Recognition (ASR) post editing as the task for applying LLM-based rewrite models and report word error rates (WER) and output length reduction rates using the span representation and the full rewrite models as the baselines.

Our contributions include the following.

- We propose a compact edit string representation with superior efficiency-accuracy trade-off than Kaneko and Okazaki (2023).
- We apply edit representation based rewriting LLMs to the task of ASR output correction. To the best of our knowledge, our work is the first to combine compact rewriting with generative LLMs to achieve substantial ASR WER reduction with manageable decoding cost.

2 Rewrite Representations

Mathematically speaking, for rewrite examples (\mathbf{x}, \mathbf{y}) , where \mathbf{x} is the input string and \mathbf{y} is the output string, there is a compression function C and an expansion function E satisfying

$$E(\mathbf{x}, C(\mathbf{x}, \mathbf{y})) = \mathbf{y} \quad (1)$$

with the constraint that $|C(\mathbf{x}, \mathbf{y})|$ has a much smaller average value than $|\mathbf{y}|$. At training time, examples are converted to $(\mathbf{x}, \hat{\mathbf{y}} = C(\mathbf{x}, \mathbf{y}))$. At inference time, the final output is obtained by applying the expansion function $E(\mathbf{x}, \mathbf{y}')$ where \mathbf{y}' is the decoding output for \mathbf{x} . The edit representations in this section differ in the choice of the function pair C and E .

2.1 Edit Span Representation

The span representation of Kaneko and Okazaki (2023) is derived from a word alignment graph \mathbf{a} between \mathbf{x} and \mathbf{y} . Given a bipartite alignment graph between the input sequence and the output sequence, we can identify pairs of word spans between the two sides. Each span pair is a local rewrite instance indicating that the source span is

substituted by the target span. In practice, the alignment is derived from the Levenshtein distance algorithm with the guarantee that the alignment links are monotonically ordered. It is always feasible to represent the entire rewrite as a sequence of local rewrite spans. For LLMs to predict the rewrites, we need a string representation of the span pairs. In their paper, the span representation is specified as $(i, j, \mathbf{y}_{\mathbf{a}(i\dots j)})$, where $\mathbf{a}(i\dots j)$ is the corresponding target span of a source span $i\dots j$ and $\mathbf{y}_{\mathbf{a}(i\dots j)}$ is the target phrase in this span. Under this representation, C is the concatenation of the ordered span representations:

$$C = \bigoplus_{(i,j) \in \mathbf{a}} (i, j, \mathbf{y}_{\mathbf{a}(i\dots j)}). \quad (2)$$

E is the program of applying the ordered local rewrites to the input sequence.

2.2 Phrase Pair Representation

The representation in Equation 2 is concise. It uses a pair of integers to represent a source span. However, this implies that LLMs have to count source tokens and generate indexing integer tokens interleaved with content tokens following the predefined format. The structured representation introduces brittleness to the model. A prediction error in the integer token sub-sequence can have a cascading effect when the entire output rewrite string is parsed and applied on the input. Instead, we resort to a natural language representation, which uses the source phrase $\mathbf{x}_{i\dots j}$ for a span (i, j) directly as the prefix for the target phrase $\mathbf{y}_{\mathbf{a}(i\dots j)}$. However, one downside of our representation is that when the span (i, j) is small, the subsequence $\mathbf{x}_{i\dots j}$ can be ambiguous, introducing errors into the expansion step. A solution is to add more context to $\mathbf{x}_{(i\dots j)}$ as well as $\mathbf{y}_{\mathbf{a}(i\dots j)}$ to make it much less likely to be ambiguous by extending the phrase pair to both the left and the right.

Formally, the new function C is

$$C = \bigoplus_{(i,j) \in \mathbf{a}} (\mathbf{W} : \mathbf{x}_{(i-k\dots j+k)}, \mathbf{y}_{\mathbf{a}(i-k\dots j+k)}), \quad (3)$$

where k is called the *dilation span*, \mathbf{W} is a natural language prompt word like `rewrite`. The expansion function E involves parsing the pattern and prefix string matching and replacement on the input.

2.3 Target Phrase Representation

The representation in Equation 3 using both source phrases and target phrases has the disadvantage

<i>source text</i>	Since we do not to bring cash to pay for the transportation fee , enormous time have been saved
<i>target text</i>	Since we do not <i>need</i> to bring cash to pay for the transportation fee , enormous time <i>has</i> been saved
<i>span</i>	<u>4 4</u> need, <u>16 17</u> has
<i>phrase pair</i>	rewrite: <u>not to</u> , not need to, rewrite: <u>time have been</u> , time has been
<i>target only</i>	rewrite: <u>not need to</u> , rewrite: <u>time has been</u>

Table 1: Example output under various representations. We underline the numerical spans or substrings to be matched against the source text in the expansion stage.

of being verbose. Dilation spans on both sides are intended to make phrases less ambiguous can exacerbate the problem. However, dilation spans on target phrases can often be sufficient for disambiguation when the span size is three or higher. We can ignore source phrases and just use dilated target phrases as they contain both anchor text in the input and replacement text in the output. The following is the new compression function.

$$C = \bigoplus_{(i,j) \in \mathbf{a}} (W : \mathbf{y}_{\mathbf{a}(i-k...j+k)}) \quad (4)$$

String matching and replacement in the implementation of function E deals with discontinuous dilation spans in the form of $\mathbf{y}_{\mathbf{a}(i-k...i-1)} \cdots \mathbf{y}_{\mathbf{a}(j+1,j+k)}$.

In Table 1, we show actual examples of edit representations. Under the target only representation, “not ... to” has two matches (“not to” and “not to bring cash to”) in the source text. For such cases, we prefer the leftmost and closest pair to break ties.

3 Experiments

We use a decoder-only LLM for the task of correcting the output of a fast first pass Automatic Speech Recognition (ASR) model. The first pass model is a streaming model that decodes as audio comes in without the full context of the future. Therefore there is enough room for error correction using a pre-trained LLM with the full ASR transcription as the input. The task can be viewed as a variant of grammatical error correction (Brockett et al., 2006).

The ASR model we use is the Google USM model (Zhang et al., 2023). The LLMs we use are the PaLM 2 Gecko and Otter models (Anil et al., 2023). We fine-tune LLMs on the LibriSpeech (Panayotov et al., 2015) training set using the dev set for hyper-parameter and checkpoint selection and the test sets for final comparisons. The ASR model is frozen in our experiments. We fine-tune the entire Transformer LLM model to minimize

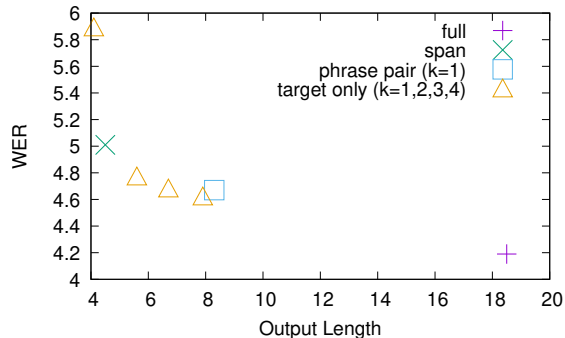
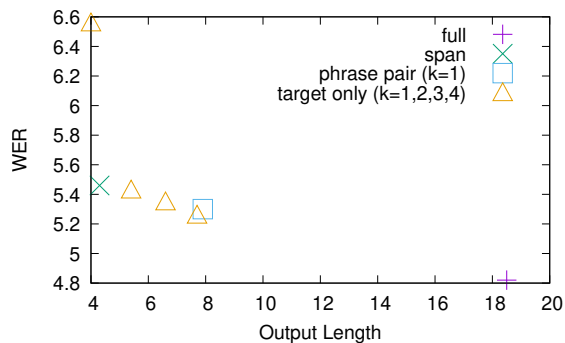


Figure 1: WER versus output length on the dev set. Top: PaLM 2 Gecko model. Bottom: PaLM 2 Otter model. *target only* with dilation span size 3 (the third \triangle from the left) is the best strategy.

the cross entropy loss on the transcription reference set given the ASR transcription generated by the frozen USM model as the prefix to the LLM decoder. We use two baselines. One is *full rewrite* model that uses the reference transcription directly. The other is *span rewrite* of Kaneko and Okazaki (2023) that uses the representation in Section 2.1. We are interested in two metrics. The quality metric is word error rate (WER) after expanding edit representations. The efficiency metric is decoder output length reduction rate.

Table 2 summarizes the main results. We show that the span representation indeed incurs more accuracy loss than the phrase representations. On the clean test set, *target only* is able to close 57% of the accuracy gap between *span* and *full*, while losing 12.5% of the length reduction rate. On the

	test-clean WER	Avg Output Length	test-other WER	Avg Output Length
<i>USM</i>	6.6	-	11.4	
<i>full</i>	2.7 (-59%)	20	6.2 (-46%)	18
<i>span</i>	3.4 (-48%)	4 (-80%)	7.5 (-34%)	5 (-72%)
<i>phrase pair</i>	3.0 (-54%)	7 (-65%)	7.1 (-38%)	10 (-44%)
<i>target only</i>	3.0 (-54%)	6 (-70%)	6.8 (-40%)	8 (-56%)

Table 2: Results of ASR (USM) post editing models based on PaLM 2 Otter. *full* has the lowest WER but has a high computational cost proportional to the average output length. *span* (Kaneko and Okazaki, 2023) is most efficient with the fewest output tokens. *target only* closes most of the WER gap between *span* and *full* while approaching the length reduction rate of *span*.

noisier other test set, *target only* is able to close 54% of the accuracy gap, while losing 22.2% of the length reduction rate.

3.1 Efficiency and Accuracy Trade-offs

In Figure 1, we plot WER versus output length for two model sizes: Gecko and Otter, and varying values of the phrase dilation hyper-parameter k in Equation 3 and Equation 4. Both the phrase pair and the target phrase only strategies yield lower WER with slightly longer outputs than the span strategy. Overall, when k is 3, the target phrase only strategy has the best trade-off. The trend stays across the two PaLM 2 model sizes.

3.2 Recovery Rate

In Section 2, we formulated the problem as selecting a pair of compression function C and expansion function E to satisfy Equation 1. The span representation is exact and unambiguous. So when E is applied, the equality is satisfied for all training examples. The phrase representations can be ambiguous and depend on the dilation spans to minimize the chance of multiple matches when the expansion function is applied. Table 3 summarizes the recovery rates, which is the percentage of examples in the dev set that satisfy Equation 1. The phrase pair representation has sufficient source context in the source phrase so that its recovery rate is very close to 100%. For the target only representation, word bigrams ($k = 2$) or trigrams ($k = 3$) surrounding target phrases are sufficient for uniquely identifying their source side counterparts in most cases.

4 Related Work

Orthogonal efforts to speed up decoding include speculative decoding (Leviathan et al., 2023; Chen et al., 2023). They leverage the overlap in output

<i>representation</i>	<i>recovery rate</i>
<i>phrase pair</i> ($k = 1$)	99.98%
<i>target only</i> ($k = 1$)	96.80%
<i>target only</i> ($k = 2$)	99.50%
<i>target only</i> ($k = 3$)	99.80%
<i>target only</i> ($k = 4$)	99.80%

Table 3: Recovery rate of phrase representations.

distributions between a less accurate faster model and a more accurate slower model as well as hardware accelerators for parallel computing. They do not incur accuracy loss and are not limited to rewriting tasks. Combining compact representations with speculative decoding has the potential for even more speedups.

LLMs have been used for ASR correction in ranking and generation (Pu et al., 2023). Leng et al. (2021) model edit operations for efficient non-autoregressive decoding. It is possible to do hybrid decoding with LLMs: predicting which spans need to be rewritten followed by auto-regressive decoding of output rewrites.

5 Conclusions

We propose two edit phrase representations for rewriting tasks that compactly represent the differences between input and output strings. We use LLMs to predict such edits and expand the edits into complete rewrites with a deterministic string matching and replacement algorithm. Our work is a further development of the span representation by Kaneko and Okazaki (2023). For the task of ASR post editing, we close 50-60% of the WER gap between the most efficient model and the most accurate model, while only slowing down decoding by 10-20% relative to the efficient representation.

295 Limitations

296 Concise edit representations presented in the paper
297 are derived from the Levenshtein distance algo-
298 rithm. The phrases are not linguistically mean-
299 ingful or optimal from machine learning point of
300 view. They are only minimal according to the edit
301 distance. Going beyond edit distance to use differ-
302 entiable functions for compression and expansion
303 is an interesting open area for research.

304 The dilation spans we use to anchor phrases in
305 the input are applied uniformly and equally on the
306 left and right of each span of interest. It is likely
307 that longer left context is more useful than right
308 context since the decoder progresses from left to
309 right.

310 We have not explored different formats for the
311 rewrite phrases.

312 The expansion stage of the target only represen-
313 tation is more involved than the other two compact
314 representations. Efficient data structures and string
315 matching algorithms are necessary to take account
316 of two discontinuous word spans.

317 Finally, we have not experimented with the latest
318 and largest LLMs. It is possible that prompt engi-
319 neering is sufficient to let these models generate
320 concise rewrites. It is to be seen if the gap between
321 full rewrite and edit representations can be reduced
322 further with very large LLMs.

323 References

324 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-
325 son, Dmitry Lepikhin, Alexandre Passos, Siamak
326 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng
327 Chen, Eric Chu, Jonathan H. Clark, Laurent El
328 Shafey, Yanping Huang, Kathy Meier-Hellstern,
329 Gaurav Mishra, Erica Moreira, Mark Omernick,
330 Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan
331 Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Her-
332 nandez Abrego, Junwhan Ahn, Jacob Austin, Paul
333 Barham, Jan Botha, James Bradbury, Siddhartha
334 Brahma, Kevin Brooks, Michele Catasta, Yong
335 Cheng, Colin Cherry, Christopher A. Choquette-
336 Choo, Aakanksha Chowdhery, Clément Crepy,
337 Shachi Dave, Mostafa Dehghani, Sunipa Dev, Ja-
338 cob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vlad
339 Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus
340 Freitag, Xavier Garcia, Sebastian Gehrmann, Lu-
341 cas Gonzalez, Guy Gur-Ari, Steven Hand, Hadi
342 Hashemi, Le Hou, Joshua Howland, Andrea Hu, Jef-
343 frey Hui, Jeremy Hurwitz, Michael Isard, Abe Itty-
344 cheriah, Matthew Jagielski, Wenhao Jia, Kathleen
345 Kenealy, Maxim Krikun, Sneha Kudugunta, Chang
346 Lan, Katherine Lee, Benjamin Lee, Eric Li, Music
347 Li, Wei Li, YaGuang Li, Jian Li, Hyeontaek Lim,

Hanzhao Lin, Zhongtao Liu, Frederick Liu, Mar-
cello Maggioni, Aroma Mahendru, Joshua Maynez,
Vedant Misra, Maysam Moussalem, Zachary Nado,
John Nham, Eric Ni, Andrew Nystrom, Alicia
Parrish, Marie Pellat, Martin Polacek, Alex Polo-
zov, Reiner Pope, Siyuan Qiao, Emily Reif, Bryan
Richter, Parker Riley, Alex Castro Ros, Aurko Roy,
Brennan Saeta, Rajkumar Samuel, Renee Shelby,
Ambrose Slone, Daniel Smilkov, David R. So,
Daniel Sohn, Simon Tokumine, Dasha Valter, Vi-
jay Vasudevan, Kiran Vodrahalli, Xuezhi Wang, Pi-
dong Wang, Zirui Wang, Tao Wang, John Wiet-
ing, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting
Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven
Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav
Petrov, and Yonghui Wu. 2023. [Palm 2 technical re-
port](#). 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364

Chris Brockett, William B. Dolan, and Michael Ga-
mon. 2006. [Correcting ESL errors using phrasal
SMT techniques](#). In *Proceedings of the 21st Interna-
tional Conference on Computational Linguistics and
44th Annual Meeting of the Association for Compu-
tational Linguistics*, pages 249–256, Sydney, Aus-
tralia. Association for Computational Linguistics. 365 366 367 368 369 370 371

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, Sandhini Agarwal, Ariel Herbert-Voss,
Gretchen Krueger, Tom Henighan, Rewon Child,
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
Clemens Winter, Christopher Hesse, Mark Chen,
Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin
Chess, Jack Clark, Christopher Berner, Sam Mc-
Candlish, Alec Radford, Ilya Sutskever, and Dario
Amodei. 2020. [Language models are few-shot learn-
ers](#). 372 373 374 375 376 377 378 379 380 381 382 383

Charlie Chen, Sebastian Borgeaud, Geoffrey Irving,
Jean-Baptiste Lespiau, Laurent Sifre, and John
Jumper. 2023. [Accelerating large language model
decoding with speculative sampling](#). 384 385 386 387

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
Maarten Bosma, Gaurav Mishra, Adam Roberts,
Paul Barham, Hyung Won Chung, Charles Sutton,
Sebastian Gehrmann, Parker Schuh, Kensen Shi,
Sasha Tsvyashchenko, Joshua Maynez, Abhishek
Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin-
odkumar Prabhakaran, Emily Reif, Nan Du, Ben
Hutchinson, Reiner Pope, James Bradbury, Jacob
Austin, Michael Isard, Guy Gur-Ari, Pengcheng
Yin, Toju Duke, Anselm Levskaya, Sanjay Ghe-
mawat, Sunipa Dev, Henryk Michalewski, Xavier
Garcia, Vedant Misra, Kevin Robinson, Liam Fe-
dus, Denny Zhou, Daphne Ippolito, David Luan,
Hyeontaek Lim, Barret Zoph, Alexander Spiridonov,
Ryan Sepassi, David Dohan, Shivani Agrawal, Mark
Omernick, Andrew M. Dai, Thanumalayan Sankara-
narayana Pillai, Marie Pellat, Aitor Lewkowycz,
Erica Moreira, Rewon Child, Oleksandr Polozov,
Katherine Lee, Zongwei Zhou, Xuezhi Wang, Bren-
nan Saeta, Mark Diaz, Orhan Firat, Michele Catasta,
Jason Wei, Kathy Meier-Hellstern, Douglas Eck, 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408

409	Jeff Dean, Slav Petrov, and Noah Fiedel. 2022.	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	468
410	Palm: Scaling language modeling with pathways.	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	469
411	Tao Fang, Shu Yang, Kaixin Lan, Derek F. Wong, Jin-	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	470
412	peng Hu, Lidia S. Chao, and Yue Zhang. 2023. Is	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	471
413	chatgpt a highly fluent grammatical error correction	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	472
414	system? a comprehensive evaluation.	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	473
415	Masahiro Kaneko and Naoaki Okazaki. 2023. Reduc-	Christina Kim, Yongjik Kim, Jan Hendrik Kirchner,	474
416	ing sequence length by predicting edit operations	Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz	475
417	with large language models.	Kondraciuk, Andrew Kondrich, Aris Konstantinidis,	476
418	Philipp Koehn, Franz Josef Och, and Daniel Marcu.	Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael	477
419	2003. Statistical phrase-based translation. In <i>Pro-</i>	Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Le-	478
420	<i>ceedings of the 2003 Conference of the North Amer-</i>	ung, Daniel Levy, Chak Ming Li, Rachel Lim,	479
421	<i>ican Chapter of the Association for Computational</i>	Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa	480
422	<i>Linguistics on Human Language Technology - Vol-</i>	Lopez, Ryan Lowe, Patricia Lue, Anna Makanju,	481
423	<i>ume 1</i> , NAACL '03, pages 48–54. Association for	Kim Malfacini, Sam Manning, Todor Markov, Yaniv	482
424	Computational Linguistics. 2003 Conference of	Markovski, Bianca Martin, Katie Mayer, Andrew	483
425	the North American Chapter of the Association	Mayne, Bob McGrew, Scott Mayer McKinney,	484
426	for Computational Linguistics on Human Language	Christine McLeavey, Paul McMillan, Jake McNeil,	485
427	Technology (HLT-NAACL 2003) ; Conference date:	David Medina, Aalok Mehta, Jacob Menick, Luke	486
428	27-05-2003 Through 01-06-2003.	Metz, Andrey Mishchenko, Pamela Mishkin, Vin-	487
429	Yichong Leng, Xu Tan, Linchen Zhu, Jin Xu, Renqian	nie Monaco, Evan Morikawa, Daniel Molling, Tong	488
430	Luo, Linqun Liu, Tao Qin, Xiangyang Li, Edward	Mu, Mira Murati, Oleg Murk, David Mély, Ashvin	489
431	Lin, and Tie-Yan Liu. 2021. Fastcorrect: Fast error	Nair, Reiichiro Nakano, Rajeev Nayak, Arvind	490
432	correction with edit alignment for automatic speech	Neelakantan, Richard Ngo, Hyeonwoo Noh, Long	491
433	recognition. <i>Advances in Neural Information Pro-</i>	Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	492
434	<i>cessing Systems</i> , 34:21708–21719.	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	493
435	Yaniv Leviathan, Matan Kalman, and Yossi Matias.	tista Parascandolo, Joel Parish, Emy Parparita, Alex	494
436	2023. Fast inference from transformers via specu-	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	495
437	lative decoding.	man, Filipe de Avila Belbute Peres, Michael Petrov,	496
438	OpenAI, Josh Achiam, Steven Adler, Sandhini Agar-	Henrique Ponde de Oliveira Pinto, Michael, Poko-	497
439	wal, Lama Ahmad, Ilge Akkaya, Florencia Leoni	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-	498
440	Aleman, Diogo Almeida, Janko Altschmidt,	ell, Alethea Power, Boris Power, Elizabeth Proehl,	499
441	Sam Altman, Shyamal Anadkat, Red Avila, Igor	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	500
442	Babuschkin, Suchir Balaji, Valerie Balcom, Paul	Cameron Raymond, Francis Real, Kendra Rim-	501
443	Baltescu, Haiming Bao, Mohammad Bavarian,	bach, Carl Ross, Bob Rotsted, Henri Roussez,	502
444	Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel	Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani	503
445	Bernadett-Shapiro, Christopher Berner, Lenny Bo-	Santurkar, Girish Sastry, Heather Schmidt, David	504
446	donoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa	Schnurr, John Schulman, Daniel Selsam, Kyla Shep-	505
447	Brakman, Greg Brockman, Tim Brooks, Miles	pard, Toki Sherbakov, Jessica Shieh, Sarah Shoker,	506
448	Brundage, Kevin Button, Trevor Cai, Rosie Camp-	Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie	507
449	bell, Andrew Cann, Brittany Carey, Chelsea Carl-	Simens, Jordan Sitkin, Katarina Slama, Ian Sohl,	508
450	son, Rory Carmichael, Brooke Chan, Che Chang,	Benjamin Sokolowsky, Yang Song, Natalie Stau-	509
451	Fotis Chantzis, Derek Chen, Sully Chen, Ruby	dacher, Felipe Petroski Such, Natalie Summers, Ilya	510
452	Chen, Jason Chen, Mark Chen, Ben Chess, Chester	Sutskever, Jie Tang, Nikolas Tezak, Madeleine B.	511
453	Cho, Casey Chu, Hyung Won Chung, Dave Cum-	Thompson, Phil Tillet, Amin Tootoonchian, Eliz-	512
454	ummings, Jeremiah Currier, Yunxing Dai, Cory De-	abeth Tseng, Preston Tuggle, Nick Turley, Jerry	513
455	careaux, Thomas Degry, Noah Deutsch, Damien	Tworek, Juan Felipe Cerón Uribe, Andrea Vallone,	514
456	Deville, Arka Dhar, David Dohan, Steve Dowling,	Arun Vijayvergiya, Chelsea Voss, Carroll Wain-	515
457	Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna	wright, Justin Jay Wang, Alvin Wang, Ben Wang,	516
458	Eloundou, David Farhi, Liam Fedus, Niko Felix,	Jonathan Ward, Jason Wei, CJ Weinmann, Ak-	517
459	Simón Posada Fishman, Juston Forte, Isabella Ful-	ila Welihinda, Peter Welinder, Jiayi Weng, Lilian	518
460	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	Weng, Matt Wiethoff, Dave Willner, Clemens Win-	519
461	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	ter, Samuel Wolrich, Hannah Wong, Lauren Work-	520
462	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao,	521
463	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Woj-	522
464	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	ciech Zaremba, Rowan Zellers, Chong Zhang, Mar-	523
465	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	vin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang	524
466	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-	525
467	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	4 technical report.	526
		Vassil Panayotov, Guoguo Chen, Daniel Povey, and	527
		Sanjeev Khudanpur. 2015. Librispeech: An asr	528
		corpus based on public domain audio books. In	529
		<i>2015 IEEE International Conference on Acoustics,</i>	530

- 531 *Speech and Signal Processing (ICASSP)*, pages
532 5206–5210.
- 533 Jie Pu, Thai-Son Nguyen, and Sebastian Stüker. 2023.
534 [Multi-stage large language model correction for](#)
535 [speech recognition](#).
- 536 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
537 Lee, Sharan Narang, Michael Matena, Yanqi
538 Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring](#)
539 [the limits of transfer learning with a unified text-to-](#)
540 [text transformer](#). *Journal of Machine Learning Re-*
541 *search*, 21(140):1–67.
- 542 Emily Reif, Daphne Ippolito, Ann Yuan, Andy Co-
543 enen, Chris Callison-Burch, and Jason Wei. 2022. [A](#)
544 [recipe for arbitrary text style transfer with large lan-](#)
545 [guage models](#). In *Proceedings of the 60th Annual*
546 *Meeting of the Association for Computational Lin-*
547 *guistics (Volume 2: Short Papers)*, pages 837–848,
548 Dublin, Ireland. Association for Computational Lin-
549 guistics.
- 550 Sascha Rothe, Jonathan Mallinson, Eric Malmi, Sebas-
551 tian Krause, and Aliaksei Severyn. 2021. [A sim-](#)
552 [ple recipe for multilingual grammatical error cor-](#)
553 [rection](#). In *Proceedings of the 59th Annual Meet-*
554 *ing of the Association for Computational Linguistics*
555 *and the 11th International Joint Conference on Nat-*
556 *ural Language Processing (Volume 2: Short Papers)*,
557 pages 702–707, Online. Association for Computa-
558 tional Linguistics.
- 559 Yu Zhang, Wei Han, James Qin, Yongqiang Wang,
560 Ankur Bapna, Zhehuai Chen, Nanxin Chen, Bo Li,
561 Vera Axelrod, Gary Wang, Zhong Meng, Ke Hu, An-
562 drew Rosenberg, Rohit Prabhavalkar, Daniel S. Park,
563 Parisa Haghani, Jason Riesa, Ginger Perng, Hagen
564 Soltau, Trevor Strohman, Bhuvana Ramabhadran,
565 Tara Sainath, Pedro Moreno, Chung-Cheng Chiu, Jo-
566 han Schalkwyk, Françoise Beaufays, and Yonghui
567 Wu. 2023. [Google usm: Scaling automatic spee-](#)
568 [ch recognition beyond 100 languages](#).