

Text Style Transfer: A Comprehensive Study on Methodologies and Evaluation

Nirali Parekh*, Siddharth Trivedi[†], and Kriti Srivastava[‡]

Department of Computer Engineering
Dwarkadas J. Sanghvi College of Engineering
Mumbai, India

*nirali25parekh@gmail.com, [†]siddharthtrivedi19@gmail.com, [‡]kriti.srivastava@djsce.ac.in

Abstract—Text Style Transfer (TST) rewords a sentence from one style (e.g. polite) to another (e.g. impolite) while conserving the meaning and content. This domain has attracted the attention of many researchers as it makes natural language generation (NLG) tasks more user-oriented. TST finds its applications widely in industry such as conversational bots and writing assistance tools. With the success of deep learning, a plethora of research works on style transfer based on Machine Learning have been proposed, developed and tested. This systematic review presents the past work on Text Style Transfer clustered into categories based on Machine Learning algorithms. It briefly explains the various subtasks within Text Style Transfer and assembles its publicly available datasets. It summarizes the evaluation metrics used for Style Transfer Tasks and throws light on practical applications and use-cases of Text Style Transfer. Finally, it expands on current challenges and points towards promising future directions for research in TST domain.

Index Terms—Text Style Transfer, Deep Learning, Natural Language Generation, Natural Language Processing, Neural Networks

I. INTRODUCTION

In the domains of Deep learning and Artificial intelligence, style transfer has recently been a hot issue of research and development. It deals with transferring the style (or attributes) of a source content into a target style. Substantial research work involving state-of-the-art algorithms [1], [2] developed for image style transfer have demonstrated astounding results. Style transfer techniques have also recently influenced the audio domain establishing methods for Music Style transfer [3]. Such works have proved the potential of deep learning techniques for style transfer in generating artificial style-transferred content. In the domain of text and linguistics, style is a highly subjective phrase that can be interchanged with the term attribute. The style-specific characteristics of the text tend to vary across situations while the style-independent content is maintained. For instance, ‘If you have any further requirements, please do not hesitate to contact me’ is used in a formal setting and can be easily paraphrased to ‘Let me know if you need anything else’ for usage in an informal context. The goal in a Text Style Transfer (TST) problem is to generate a style-controlled text in a target style while preserving the semantics and content of the source text. Text style transfer models can be formulated [4] as: $p(x | a, x')$, where x' is a source text with attribute value a' and x is the style-transferred text with target attribute a .

TST methods have developed from classic replacement and template-based approaches to neural network-based strategies as deep-learning progresses. TST can also be formulated as a Natural Language Generation (NLG) problem as it extend NLG techniques while manipulating the attributes of the text. A wide range of deep-learning techniques employed for TST tasks like Adversarial Learning, Sequence-to-Sequence Learning and Reinforcement learning-based methods which are covered in detail in section III. Experimentation on Text Style Transfer techniques is largely classified based on subtasks under the domain. List of some common subtasks in TST is elucidated in the Table I.

TABLE I: Subtasks within Text Style Transfer

Subtask	Attribute transfer	Example
Formality	Formal	Please accept our apologies for any inconvenience
	Informal	We’re sorry for the trouble
Sentiment	Positive	I admire my college professors a lot
	Negative	I hate my college professors
Politeness	Polite	Sorry, I’m a bit busy right now
	Impolite	Leave me alone
Simplicity	Complicated	Can I acquire assistance in deciphering this conundrum?
	Simple	Can you help me solve this problem?
Gender	Masculine	My wife went to the mall to buy a skirt
	Feminine	My husband went to the mall to buy a shirt
Authorship	Shakespearean	Hast thou slain tybalt?
	Modern	Have you killed tybalt?

Impactful applications of Text Style Transfer in NLP research, like paraphrasing, as well as commercial uses such as AI-assisted writing tools have driven the burgeoning interest of NLP researchers into this domain. This rapid growth in TST research has produced a variety of datasets, implementation

algorithms and evaluation metrics for this task, but also at the same time lacks a sense of standardization for the same. The objective of this review paper is to present an account of various corpora and TST methodologies that could facilitate further research and uniformity in the field of TST. The contributions of this work are:

- 1) We conduct a comprehensive survey that reviews recent works on TST based on machine learning.
- 2) We describe the various machine learning architectures and evaluation metrics used in Text Style Transfer.
- 3) We outline some of the academic and industrial applications of TST and point towards possible future directions for research.

The organization of this paper starts with a discussion on some of the publicly available datasets is done in section II. The Machine Learning algorithms used in TST research are explored in section III. section IV presents few metrics used for evaluating TST algorithms and section V outlines some applications of TST. Finally, we conclude in section VI and discuss some open issues and research scope in TST.

II. DATASETS

A. Parallel and Non-parallel data

The datasets available for Text Style Transfer are classified into two categories widely based on data used for training.

1) *Parallel*: In parallel data, the texts of both the source and target style are available. Here, simple machine translation techniques such as sequence-to-sequence can be used. The issue with parallel data is that they are not readily available in various sub-styles and data collection is very expensive.

2) *Non-parallel*: Costly and scarce parallel data has led to researchers to utilize non-parallel data for Text Style Transfer. Non-parallel corpus consists of non-matching texts from source and target styles. This category of data is readily available in various styles and hence a large proportion of works on TST utilizes non-parallel datasets.

B. Available benchmark Datasets for TST

Table II records various publicly available datasets distinguished by their subtasks, sizes, whether they are parallel or non-parallel, and annotation method incase they are parallel.

III. METHODOLOGIES

A. Parallel

Style transfer by means of a style parallel corpus is considered a monolingual machine translation task.

1) *Sequence To Sequence*: A seq2seq model known as encoder-decoder architecture converts sequences from one domain to another. Basis to many machine translation algorithms, many research works on TST utilized seq2seq neural networks on parallel datasets. Hence, a seq2seq model is trained such that the encoder is input is the text of source style, and output is the corresponding text of target style. RNN layers act as encoder process the input sentence of source style and recover the state to serve as context for the decoder. For the next step,

TABLE II: Publicly available datasets for TST

Subtask	Dataset	Size	Parallel	Domain	Annotation
Formality	GYAFC	52K	✓	Yahoo Answers (online)	Manual
Politeness	Politeness	1.39M	✗	Emails	-
Gender	Yelp	2.5M	✗	Restaurant Reviews	-
Humor & Romance	Flickr style	5K	✓	image captions	Manual
biasedness	Wiki Neutrality	181K	✓	Wikipedia (online)	Automatic
Toxicity	Twitter	58K	✗	tweets (online)	-
Toxicity	Reddit	224K	✗	politics threads	-
Authorship	Shakespeare	18K	✓	literature, SparkNotes	Automatic
Sentiment	Yelp	150K	✗	restaurant reviews	-
sentiment	Amazon	277K	✗	clothing reviews	-
Fluency	SWBD	192K	✓	telephonic conversations	Manual
Politics	Political	540K	✓	facebook posts	-

another layer of RNNs acting as decoder predicts the output sentence of the target style.

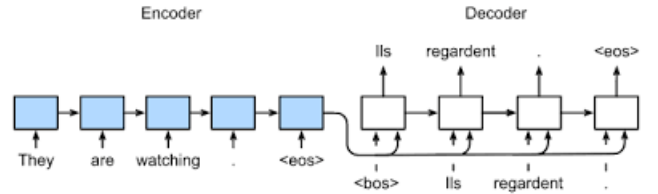


Fig. 1: Sequence-to-sequence attention-based model

In their work, Wei Xu [10] present some early work on the task of rephrasing text in a particular style. H. Jhamtani et al. [9] employ a seq2seq neural network to convert Modern English to Shakespearean text. Using dictionary to map modern and shakespearean English, they utilized a basic encoder-decoder architecture. Carlson et al. [11] utilized a seq2seq model with attention mechanism for converting a prose to Bible style text. Some later works have experimented techniques like data augmentation along with the sequence-to-sequence models [20], [21]. Nikolav et al. [21] work on the task of text simplification from regular Wikipedia to simpler version. They collect two pseudo-parallel corpus pairs from sources like technical articles and news websites and propose

TABLE III: Summary of some previous works in TST - their methodologies, data sources, and evaluation practices

Ref	ML Method	Contribution	Subtask	Data	Automatic Evaluation				Human Evaluation
					Overall	Style Transfer Accuracy	Content Preservation	Fluency	
[5]	Adversarial Learning	Two models- multi-decoder and style embedding	paper- news title, sentiment transfer	He and McAuley	-	Lstm sigmoid classifier	Cosine distance between source and target sentence embeddings	-	✓
[6]	Adversarial Learning	incorporates auxiliary and adversarial objectives	sentiment transfer	Yelp reviews, Amazon reviews	Geometric Mean of STA, WO and 1/PPL	CNN classifier	Cosine similarity, Unigram word overlap	Perplexity by trigram Kneser-Ney language model	✓
[7]	Adversarial Learning	word-level conditional architecture and a two-phase training procedure	sentiment transfer, tense transfer	Yelp Reviews, Yelp Tense, Amazon Reviews	-	CNN classifier	BLEU	Perplexity by bi-directional Lstm	✓
[8]	Adversarial Learning	adversarially regularized autoencoders (ARAE)	sentiment transfer, topic transfer	SNLI corpus, Yahoo dataset, Yelp Reviews	-	fastText classifier	BLEU	Perplexity	✓
[9]	Sequence-to-Sequence	external dictionaries that maps Shakespearean to modern English words	Old-Modern English	Shakespeare dataset	BLEU, and PINC	-	-	-	✗
[10]	Sequence-to-Sequence	translation model with a language model	Paraphrasing	Shakespeare dataset	-	Cosine similarity, Language model, Logistic regression	BLEU	-	✓
[11]	Sequence-to-Sequence	Encoder-decoder recurrent neural networks	Old-Modern English	Various English Bible versions	BLEU, and PINC	-	-	-	✗
[12]	Keyword Replacement	Delete, Retrieve, Generate	Normal- Romantic, Sentiment transfer	Yelp Review, Amazon Review, Captions dataset	-	LSTM- based classifier	BLEU	-	✓
[13]	Keyword Replacement	Transformer that leverages DRG framework	sentiment transfer, gender transfer, political slant	Yelp Reviews, Amazon, Reviews, Captions dataset, Political dataset, Gender dataset	GLEU	FastText style classifier	BLEU	Perplexity by GPT-2	✓

an unsupervised technique LHA for extracting pseudo parallel sentences from the sources. However, seq2seq approach requires parallel corpus, and is hence challenging due to its scarcity.

B. Non-Parallel

1) *Keyword Replacement*: Certain keywords in texts often are indicative of sentiments and tone underneath. For instance, words like “wow”, “excellent” have a positive meaning and “awful” and “worst” have a negative connotation. Hence, some works [12], [13] introduce the use of machine learning models to replace these words using Natural language generative models.

An initial work of the keyword replacement approach is the Delete-Retrieve-Generate framework [12]. Firstly, the model identifies all the style-connotated words like “bad”, “rude” etc.

from the source sentences (eg. negative sentiment). Then it eliminates these words and only the content indicative words are left, such as “hotel” or “shirt”. Next, reference texts similar to the content-related words are retrieved from the target sentiment corpus (here, eg. positive). Then, the model similarly gets all the style-attributed words and combines them with the content words extracted previously. This combining is done with seq2seq models

The authors of [22] uses attention mechanisms with Deep Learning architectures for Natural Language Generation tasks. Recent works explored hybrid architectures using keyword replacement methods with cycled reinforcement losses to iteratively transfer the style of the text while maintaining the content. Since the words replaced can be visually inspected by a human, it imparts explainability to the models. The parts of text modified can be examined to understand the performance

TABLE IV: Summary of some previous works in TST - their methodologies, data sources, and evaluation practices (Contd.)

Ref	ML Method	Contribution	Subtask	Data	Automatic Evaluation				Human Evaluation
					Overall	Style Transfer Accuracy	Content Preservation	Fluency	
[14]	Keyword Replacement, Reinforcement Learning	cycled reinforcement learning method	sentiment transfer	Amazon Reviews Yelp Reviews	G-score: Geometric Mean of ACC and BLEU	CNN classifier	BLEU	-	✓
[15]	Unsupervised Learning	encoder-decoder architecture reinforced through auxiliary modules called scorers	formality transfer	emails, english prose essays	-	Encoder-based neural classifier	Cosine similarity	Perplexity by 4-gram back-off model using KenLM	✓
[16]	Back-translation	personalized SMT models in automatic translation	gender transfer	TED talks transcripts Euro-parliament corpus	-	SVM classifier	-	-	✓
[17]	Adversarial Learning, Keyword Replacement	Focused on 3 models: CAEE, ARAE, DAR models	sentiment transfer	Yelp Reviews	-	Earth Mover's Distance	BLEU, METEOR, Word Mover's distance	Neural logistic regression classifiers	✓
[18]	Adversarial Learning	tag-and-generate pipeline - tagger extracts content and generator converts style	politeness transfer, Normal - Romantic, sentiment, gender transfer, Political slant	Enron Email, Captions, Yelp reviews, Amazon reviews	-	LSTM classifier	BLEU, METEOR, ROUGE	-	✓
[19]	Unsupervised Learning, Back-translation	TGLS (Text Generation by Learning from Search), framework	paraphrase generation, formality transfer	GYAFC	-	classifier based on RoBERTa features	BLEU, iBLEU	Perplexity by GPT-2	✓

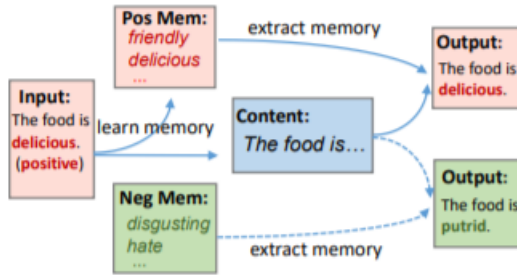


Fig. 2: Overview of the model proposed by [22] using self-attention and keyword replacement

of such techniques.

2) *Adversarial Learning*: Another effective method for TST is adversarial learning which separates the text's style and content data for transferring the text to target style. An early work proposed by Fu et al. [5] uses an adversarial framework of two models.

The first model consists of a single encoder and multiple decoders. The style of text input is learned by the encoder representations which in turn trains multiple decoders to decipher the representation and output the text in target style. In the second model, encoder behaves similar to the first model, and the decoder outputs the target style by concatenating the

embedding of the encoder to parameter representation.

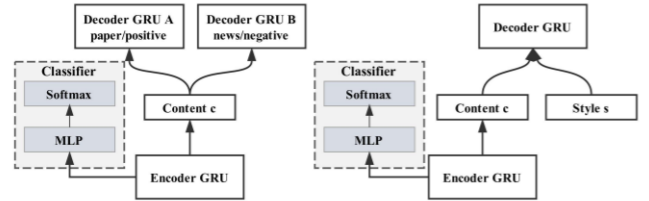


Fig. 3: Two adversarial learning based models proposed by [5]

Similar to the previous methods, a variety of hybrid architectures have been proposed for adversarial-based frameworks. Recently, works on the usage of autoencoders for TST models implementing adversarial learning to refine the performance of style transfer tasks is studied. For example, some works [6]–[8], [23] implemented cycled-consistency loss where the generated text is again fed to the model and the output sentence is compared to the original input text. This way, the loss is cycled or transferred to the model to generate more accurate transferred sentences.

3) *Reinforcement Learning*: Reinforcement learning works on the idea that reward functions guide the decision of deep learning models instead of loss functions. The parameters of a reinforcement-based model change in a way that the estimated reward of the output style-transferred text is maximized. An

attention-based model proposed by [24] uses encoder-decoder architecture to perform Text Style Transfer. Its proposed generator and evaluator based method is shown in Figure 4. Here, three rewards to guide the output text of the desired style are introduced. Style classifiers, semantic and language models are employed to impart style, semantic and language feedback to the model respectively.

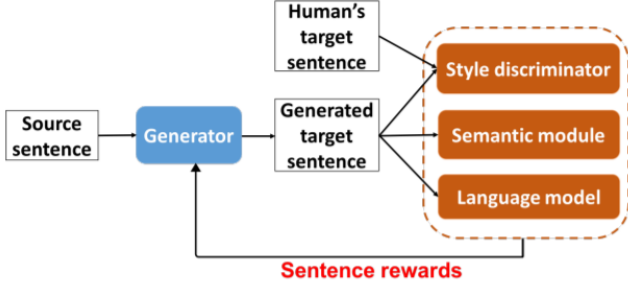


Fig. 4: Overview of the model proposed by [24] where rewards are returned to the generator

In [25], the authors propose a dual reinforcement learning framework where one seq2seq model learns source style to target style embeddings. The second seq2seq learns the target to source style parameters. This dual task is designed to provide style and reconstruction rewards and their average is used to provide feedback to the model. Thus, reinforcement learning without any parallel corpus is used for style transfer tasks.

4) *Backtranslation*: Back-translation means translating a text of target style to the source style and mixing both original source and back-translated text to train the model. The work by authors of [16] researched the factors like gender that are often obscured in tasks like Machine Language Translation. The back-translation method used for Text Style Transfer by [26] is shown in Figure 5. Their approach was to first rephrase the sentences and retain only the content words thereby eliminating the style information. Using an NMT model, they translated English text into another language and then back into the English language. It is proposed with the understanding that the stylistic properties are lost during the back and forth translation while preserving the content. In [19], the authors propose a for text generation framework with two stages. First stage of Learning simulated annealing search to generate pseudo-input sentences and a generator undergoes training via supervised learning. In the next stage, iterations are performed of beam search process for model improvement.

The authors in [27] employ a two-fold framework where first a pseudo-parallel dataset is formed using word embeddings and latent representations similarities. Thus this unsupervised task of style transfer is now employed to back translation-based deep learning models. A style classifier is utilized to enhance the performance of the model. Like the methods using adversarial models, back-translation models also proves less efficient in style-content agreement.

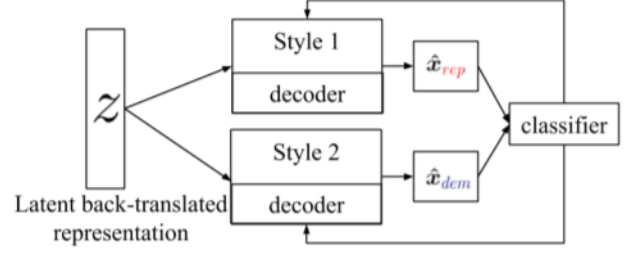


Fig. 5: Back-translation method and style classifier used by [26]

C. Unsupervised methods

The previous models using parallel or non-parallel data are proposed for supervised settings for style transfer tasks. Recent works have also explored Machine Learning techniques in purely unsupervised settings where no labelled style-text corpus is provided to the model. There are relatively fewer works proposed to perform TST in such a way [15], [28], [29].

Initial study conducted by Radford et al. [28] made use of the characteristics RNNs and LSTMs wherein training of such models is done on the UTF8 byte of the input text. This preprocessing allowed the researchers to identify and modify the neuron-level embeddings of the Deep Learning models. Text formalization performed by Jain et al. [15] uses an unsupervised method with unlabelled corpus. They make use of external language processing tools called scorers to provide style information. The information learned by the encoder-decoder is backpropagated to the model for computing the loss. The output scores determine the formality level of the text and the resulting output text. Shen et al. in their work [29] proposed adversarial autoencoders called AAE with denoising models for mapping of similar latent representations. These models perform sentiment transfer by computing a vector using these representations. There is still much research potential in natural language generation based unsupervised methods that can be extrapolated to other TST subtasks.

IV. EVALUATION TECHNIQUES

Measuring the true efficacy of the TST models is one of the most challenging tasks in the domain of text style transfer. At present, there are no standard automatic evaluation metrics that are being followed conventionally and still no metric exists that can outperform human evaluation. Since human evaluation can be onerous and expensive, there is a compelling need for appropriate automatic evaluation practices. Various evaluation metrics for TST have been proposed previously which mainly focus on three aspects to measure the effectiveness of a TST model:

- 1) *Style transfer strength*: the ability to convert a source style to the desired target style.

- 2) Content preservation: the extent to which the original content is preserved.
- 3) Fluency (naturalness): the ability to generate fluent sentences.

A TST model should perform deftly in all these three criteria of evaluation. Underperformance in any of the three aspects would rather deem the model ineffective. For instance, if the algorithm successfully transfers a positive sentiment sentence, “The teachers in this university are excellent” to a negative sentiment sentence with fluent grammar, “The waiters in the restaurant were very rude”, the algorithm, here, doesn’t preserve the original meaning and hence is considered to be inadequate.

A. Automatic Evaluation

1) *Style transfer strength*: This criterion of evaluation deals with how well the style of a given text is transferred to the target style. In most past works, transfer strength is tested by making use of pre-trained classifiers. A lot of previous papers [20], [25], [30] have used TextCNN, a sentence-level text classifier trained over pre-trained word vectors, proposed by [31]. An LSTM classifier, first used for this task in [5], is employed to measure the transfer strength in [24], [32]. Some works [33], [34] have made use of fastText [35], which shows identical performance as the deep learning methods, while being faster in speed, for the style classification task. Another alternate metric proposed in [17] for measuring transfer strength is to compute the Earth Mover’s Distance [36] between the source text and the style transferred output. This metric can be used to handle even non-binary classification and exhibited a higher correlation with human evaluation than fastText and TextCNN classifiers.

2) *Content Preservation*: The metrics under this criterion of evaluation are concerned with the measurement of the extent of original content that is preserved after style transfer. The most widely used metric to evaluate content similarity is the BLEU score [37], originally proposed for the evaluation of machine translation tasks. [14], [38] have computed the BLEU score between the source text and the transferred text to measure content preservation in transferred sentences. In [5], the authors have calculated the cosine similarity between the source and the transfer sentence embeddings by leveraging the pre-trained word embeddings by GloVe [39]. Another popular metric for this specific task is the METEOR score [40] used in [18], [41]. [17] proposes to practice style masking or style removal i.e., to remove or mask the style attribute words from the source and transferred sentences, before using the content preservation metrics. The authors, in [42] have done an extensive comparison of fourteen content similarity metrics over two style transfer datasets and suggested that Word Mover’s Distance (WMD) [43] and L2 distance based on ELMo [44] are the best performing metrics for measuring content preservation in style transfer tasks.

3) *Fluency (naturalness)*: For any natural language generation model to be efficient, it must possess the ability to produce fluent human-like text. Most commonly, like in [6],

[13], [15], a language model is trained and employed to compute perplexity score (PPL), where lower scores indicate higher fluency. Although researchers have relied on using the perplexity scores for evaluating fluency of the TST models in the past, [17] showed that perplexity scores exhibited a very low correlation with human-evaluated scores and instead adversarial neural classifiers should be employed for the task of evaluating naturalness.

B. Human Evaluation

A lot of previous works [5], [24], [27], [30] have incorporated the use of human evaluation along with the automatic evaluation metrics discussed above. Generally, a common procedure is followed where the evaluators are asked to rate randomly selected style transferred outputs. The scale of rating varies across different works but essentially the approach is similar where higher scores indicate better performance. The human raters can provide overall scores or give scores separately based on the three criteria i.e., style transfer strength, content preservation, and fluency. Human evaluation, apart from being expensive and cumbersome, for the style transfer task can often be subjective in nature depending on how the rater interprets the styles under consideration. Therefore, it is also often not comparable across different methods due to this subjectiveness. Notwithstanding the cons, human evaluation is extremely important to be carried out in the contemporary scenario of the TST domain along with the automatic evaluation, to study the correlation between various automatic metrics and human judgment. It will play an important role in the evolution of automatic metrics for TST tasks and set benchmarks for future comparison.

V. APPLICATIONS

The TST domain holds a number of research-based and consumer-based applications. TST techniques can be applied to aid research in the Natural language domain as well as be implemented to develop user-assisting products for applications concerned with conversation, literature, and languages. This section summarizes the potential applications of the TST techniques.

1) *Text Augmentation*: Text augmentation is the generation of text similar to the existing textual data in order to expand the training corpus when there is a lack of data. Some deep learning-based works [45], [46] do exist for text augmentation but TST as a text augementer tool has not been explored much. TST techniques are capable of being employed for this task as they can produce text in different styles while preserving the content.

2) *Paraphrasing and Summarization*: Paraphrasing a text or a passage means to restate the information in the text in other words while preserving its meaning. This task is quite synonymous with TST, where the text is generated in a different style while the content is preserved. [42] explores the similarity between paraphrasing and style transfer by comparing the evaluation metrics for content preservation. Stylistic summarization is also one of the applications of

TST. [47] has proposed a method to generate style-specific headlines for a text which produces summarized headlines in three different styles for the same text.

3) *Conversational AI*: With the advancements in Natural language processing and Artificial intelligence, there have been rapid developments in AI-powered chatbots and conversational agents lately. Past work in this area [48] has pointed out how the conversational style in chatbots can affect user interaction. TST techniques can be applied to these chatbots to generate persona-based responses [49] and make them flexible for conversation in various styles.

4) *Text Anonymization*: This area deals with protecting the privacy of a user by preventing leaks of private-attribute information related to the user. It has become essential in the contemporary world to secure sensitive textual data from malicious attackers. TST techniques can be encountered as a potential solution in this area for manipulating the text and obscuring the user identity [32], [50].

5) *Writing assistance tools*: TST techniques can also be extensively used for the development of various computer-aided writing applications used for assisting users to write with precision. TST methods can detect the formality or politeness of a text and help in guiding the author to modify the text as per the required style for particular use-cases, for instance, a business agreement. [51] has worked upon technique for simplification of texts which can help certain users to comprehend complex passages with some ease.

VI. DISCUSSION AND CONCLUSION

Research in the field of TST is challenged by a few obstacles such as scarcity of publicly available, benchmark parallel corpora and the absence of standardized automatic evaluation metrics. Currently, there are no automatic evaluation methods developed that can surpass human evaluation, which proves to be expensive and cumbersome. Hence, a prospective research area in TST is using unsupervised machine learning methods to surpass the issue of the absence of large parallel data. Also, while a significant portion of the work is applicable to English corpora, TST's potential to extrapolate it to other languages should not be neglected. In this paper, we have provided a comprehensive review of the existing literature and numerous machine learning methods employed for the task. The review also discussed various benchmark datasets, evaluation metrics used for TST and presented an overview of several important applications where TST methods are of use. This paper can serve as a reference for NLP researchers and is aimed to provide an all-inclusive understanding of TST to facilitate and promote further research in this field.

REFERENCES

- [1] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2242–2251.
- [2] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [3] O. Cifka, U. Şimşekli, and G. Richard, "Groove2groove: One-shot music style transfer with supervision from synthetic data," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2638–2650, 2020.
- [4] D. Jin, Z. Jin, Z. Hu, O. Vechtomova, and R. Mihalcea, "Deep learning for text style transfer: A survey," *CoRR*, vol. abs/2011.00416, 2020. [Online]. Available: <https://arxiv.org/abs/2011.00416>
- [5] Z. Fu, X. Tan, N. Peng, D. Zhao, and R. Yan, "Style transfer in text: Exploration and evaluation," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, Apr. 2018. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/11330>
- [6] V. John, L. Mou, H. Bahuleyan, and O. Vechtomova, "Disentangled representation learning for non-parallel text style transfer," *arXiv preprint arXiv:1808.04339*, 2018.
- [7] C.-T. Lai, Y.-T. Hong, H.-Y. Chen, C.-J. Lu, and S.-D. Lin, "Multiple text style transfer by using word-level conditional generative adversarial network with two-phase training," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 3579–3584.
- [8] J. Zhao, Y. Kim, K. Zhang, A. Rush, and Y. LeCun, "Adversarially regularized autoencoders," in *International conference on machine learning*. PMLR, 2018, pp. 5902–5911.
- [9] H. Jhamtani, V. Gangal, E. Hovy, and E. Nyberg, "Shakespeareizing modern language using copy-enriched sequence-to-sequence models," *arXiv preprint arXiv:1707.01161*, 2017.
- [10] W. Xu, A. Ritter, W. B. Dolan, R. Grishman, and C. Cherry, "Paraphrasing for style," in *Proceedings of COLING 2012*, 2012, pp. 2899–2914.
- [11] K. Carlson, A. Riddell, and D. Rockmore, "Evaluating prose style transfer with the bible," *Royal Society open science*, vol. 5, no. 10, p. 171920, 2018.
- [12] J. Li, R. Jia, H. He, and P. Liang, "Delete, retrieve, generate: A simple approach to sentiment and style transfer," *arXiv preprint arXiv:1804.06437*, 2018.
- [13] A. Sudhakar, B. Upadhyay, and A. Maheswaran, "Transforming delete, retrieve, generate approach for controlled text style transfer," *arXiv preprint arXiv:1908.09368*, 2019.
- [14] J. Xu, X. Sun, Q. Zeng, X. Zhang, X. Ren, H. Wang, and W. Li, "Unpaired sentiment-to-sentiment translation: A cycled reinforcement learning approach," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2018. [Online]. Available: <https://doi.org/10.18653/v1/p18-1090>
- [15] P. Jain, A. Mishra, A. P. Azad, and K. Sankaranarayanan, "Unsupervised controllable text formalization," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 6554–6561.
- [16] E. Rabinovich, S. Mirkin, R. N. Patel, L. Specia, and S. Wintner, "Personalized machine translation: Preserving original author traits," *arXiv preprint arXiv:1610.05461*, 2016.
- [17] R. Mir, B. Felbo, N. Obradovich, and I. Rahwan, "Evaluating style transfer for text," in *Proceedings of the 2019 Conference of the North Association for Computational Linguistics*, 2019. [Online]. Available: <https://doi.org/10.18653/v1/n19-1049>
- [18] A. Madaan, A. Setlur, T. Parekh, B. Poczos, G. Neubig, Y. Yang, R. Salakhutdinov, A. W. Black, and S. Prabhunoye, "Politeness transfer: A tag and generate approach," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.169>
- [19] J. Li, Z. Li, L. Mou, X. Jiang, M. R. Lyu, and I. King, "Unsupervised text generation by learning from search," *arXiv preprint arXiv:2007.08557*, 2020.
- [20] Z. Jin, D. Jin, J. Mueller, N. Matthews, and E. Santus, "IMaT: Unsupervised text attribute transfer via iterative matching and translation," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, 2019. [Online]. Available: <https://doi.org/10.18653/v1/d19-1306>
- [21] N. I. Nikolov and R. H. Hahnloser, "Large-scale hierarchical alignment for data-driven text rewriting," *arXiv preprint arXiv:1810.08237*, 2018.
- [22] Y. Zhang, J. Xu, P. Yang, and X. Sun, "Learning sentiment memories for sentiment modification without parallel data," *arXiv preprint arXiv:1808.07311*, 2018.

- [23] L. Chen, S. Dai, C. Tao, D. Shen, Z. Gan, H. Zhang, Y. Zhang, and L. Carin, "Adversarial text generation via feature-mover's distance," *arXiv preprint arXiv:1809.06297*, 2018.
- [24] H. Gong, S. Bhat, L. Wu, J. Xiong, and W. mei Hwu, "Reinforcement learning based text style transfer without parallel training corpus," in *Proceedings of the 2019 Conference of the North Association for Computational Linguistics*, 2019. [Online]. Available: <https://doi.org/10.18653/v1/n19-1320>
- [25] F. Luo, P. Li, J. Zhou, P. Yang, B. Chang, X. Sun, and Z. Sui, "A dual reinforcement learning framework for unsupervised text style transfer," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, Aug. 2019. [Online]. Available: <https://doi.org/10.24963/ijcai.2019/711>
- [26] S. Prabhunoye, Y. Tsvetkov, R. Salakhutdinov, and A. W. Black, "Style transfer through back-translation," *arXiv preprint arXiv:1804.09000*, 2018.
- [27] Z. Zhang, S. Ren, S. Liu, J. Wang, P. Chen, M. Li, M. Zhou, and E. Chen, "Style transfer as unsupervised machine translation," *ArXiv*, vol. abs/1808.07894, 2018.
- [28] A. Radford, R. Jozefowicz, and I. Sutskever, "Learning to generate reviews and discovering sentiment," *arXiv preprint arXiv:1704.01444*, 2017.
- [29] T. Shen, J. Mueller, R. Barzilay, and T. Jaakkola, "Educating text autoencoders: Latent representation guidance via denoising," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8719–8729.
- [30] T. Shen, T. Lei, R. Barzilay, and T. Jaakkola, "Style transfer from non-parallel text by cross-alignment," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 6833–6844.
- [31] Y. Kim, "Convolutional neural networks for sentence classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2014. [Online]. Available: <https://doi.org/10.3115/v1/d14-1181>
- [32] T. Gröndahl and N. Asokan, "Effective writing style imitation via combinatorial paraphrasing," *CoRR*, vol. abs/1905.13464, 2019. [Online]. Available: <http://arxiv.org/abs/1905.13464>
- [33] Y. Liu, G. Neubig, and J. Wieting, "On learning text style transfer with direct rewards," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, 2021. [Online]. Available: <https://doi.org/10.18653/v1/2021.naacl-main.337>
- [34] N. Dai, J. Liang, X. Qiu, and X. Huang, "Style transformer: Unpaired text style transfer without disentangled latent representation," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019. [Online]. Available: <https://doi.org/10.18653/v1/p19-1601>
- [35] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, "Bag of tricks for efficient text classification," in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*. Valencia, Spain: Association for Computational Linguistics, Apr. 2017, pp. 427–431. [Online]. Available: <https://aclanthology.org/E17-2068>
- [36] Y. Rubner, C. Tomasi, and L. J. Guibas, "A metric for distributions with applications to image databases," in *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*. IEEE, 1998, pp. 59–66.
- [37] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "BLEU," in *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics - ACL '02*. Association for Computational Linguistics, 2001. [Online]. Available: <https://doi.org/10.3115/1073083.1073135>
- [38] M. Shang, P. Li, Z. Fu, L. Bing, D. Zhao, S. Shi, and R. Yan, "Semi-supervised text style transfer: Cross projection in latent space," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, 2019. [Online]. Available: <https://doi.org/10.18653/v1/d19-1499>
- [39] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, 2014. [Online]. Available: <https://doi.org/10.3115/v1/d14-1162>
- [40] S. Banerjee and A. Lavie, "METEOR: An automatic metric for MT evaluation with improved correlation with human judgments," in *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*. Ann Arbor, Michigan: Association for Computational Linguistics, Jun. 2005, pp. 65–72. [Online]. Available: <https://aclanthology.org/W05-0909>
- [41] R. Shetty, B. Schiele, and M. Fritz, "A4nt: Author attribute anonymity by adversarial training of neural machine translation," in *27th USENIX Security Symposium (USENIX Security 18)*. Baltimore, MD: USENIX Association, Aug. 2018, pp. 1633–1650. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity18/presentation/shetty>
- [42] I. P. Yamshchikov, V. Shibaev, N. Khlebnikov, and A. Tikhonov, "Style-transfer and paraphrase: Looking for a sensible semantic similarity metric," *ArXiv*, vol. abs/2004.05001, 2021.
- [43] M. Kusner, Y. Sun, N. Kolkin, and K. Weinberger, "From word embeddings to document distances," in *Proceedings of the 32nd International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, F. Bach and D. Blei, Eds., vol. 37. Lille, France: PMLR, 07–09 Jul 2015, pp. 957–966. [Online]. Available: <https://proceedings.mlr.press/v37/kusnerb15.html>
- [44] M. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, "Deep contextualized word representations," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, 2018. [Online]. Available: <https://doi.org/10.18653/v1/n18-1202>
- [45] A. Amin-Nejad, J. Ive, and S. Velupillai, "Exploring transformer text generation for medical dataset augmentation," in *Proceedings of the 12th Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, May 2020, pp. 4699–4708. [Online]. Available: <https://aclanthology.org/2020.lrec-1.578>
- [46] V. Atliha and D. Šešok, "Text augmentation using BERT for image captioning," *Applied Sciences*, vol. 10, no. 17, p. 5978, Aug. 2020. [Online]. Available: <https://doi.org/10.3390/app10175978>
- [47] D. Jin, Z. Jin, J. T. Zhou, L. Orii, and P. Szolovits, "Hooks in the headline: Learning to generate headlines with controlled styles," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.456>
- [48] C. Liebrecht, L. Sander, and C. van Hooijdonk, *Too Informal?: How a Chatbot's Communication Style Affects Brand Attitude and Quality of Interaction*. Springer, 2020, vol. 12604, pp. 16–31.
- [49] J. Li, M. Galley, C. Brockett, G. Spithourakis, J. Gao, and B. Dolan, "A persona-based neural conversation model," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2016. [Online]. Available: <https://doi.org/10.18653/v1/p16-1094>
- [50] S. Reddy and K. Knight, "Obfuscating gender in social media writing," in *Proceedings of the First Workshop on NLP and Computational Social Science*. Association for Computational Linguistics, 2016. [Online]. Available: <https://doi.org/10.18653/v1/w16-5603>
- [51] Y. Cao, R. Shui, L. Pan, M.-Y. Kan, Z. Liu, and T.-S. Chua, "Expertise style transfer: A new task towards better communication between experts and laymen," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020. [Online]. Available: <https://doi.org/10.18653/v1/2020.acl-main.100>