

Style-Specific Neurons for Steering LLMs in Text Style Transfer

Wen Lai^{1,2}, Viktor Hangya³, Alexander Fraser^{1,2}

¹School of Computation, Information and Technology, Technical University of Munich, Germany

²Munich Center for Machine Learning, Germany

³Center for Information and Language Processing, LMU Munich, Germany

12th November, 2024



1 Introduction

2 Related Work

3 Method

4 Experiments

5 Results

6 Analysis

7 Conclusion

Background (I)

- Text style transfer (TST) aims to transform text from a **source style** to a **target style** while **maintaining the original content** and **ensuring the fluency** of the generated text.

Style Features	Example	
Formality	Formal:	Would you like to get a drink?
	Informal:	Wanna get a drink?
Politeness	Polite:	I see your point, but I have a different perspective.
	Impolite:	You're wrong.
Sentiment	Positive:	I'd enjoy something more exciting.
	Negative:	This is boring.

Background (II)

- LLMs tends to directly copy a significant portion of the input text to the output without effectively changing its style.

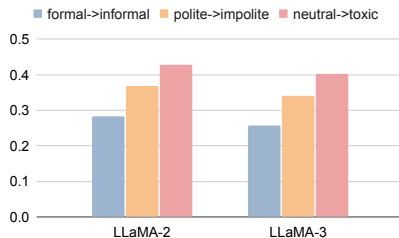


Fig. Copy ratio for LLMs (e.g., LLaMA-2 and LLaMA-3) on TST task.

Enhancing the generation of words that align with the target style during the decoding process remains a significant challenge in TST.

Background (III)

- Neuron analysis aims to identify and understand the roles of individual neurons within a neural network.

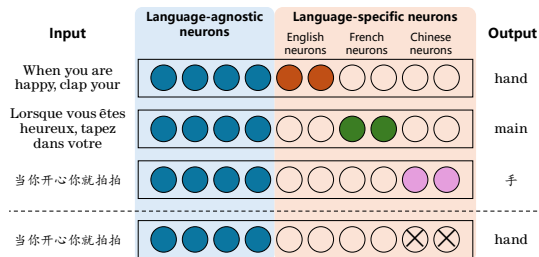


Figure directly taken from Tang et al., 2024

- LLMs exist language-specific neurons that can process the vocabulary, grammar, and idioms of a particular language.
- Deactivating the language-specific neurons leads to a remarkable degradation in the model's understanding and generation abilities.

Research Question

- **Q1:** Do LLMs possess neurons that specialize in processing style-specific text?
- **Q2:** If such neurons exist, how can we optimize their utilization during the decoding process to steer LLMs in generating text that faithfully adheres to the target style?

Text Style Transfer

- Traditional TST task: parallel data, non-parallel data ([Jin et al., 2022](#))
- Recent work on using LLMs
 - Fine-tuning ([Mukherjee et al., 2024](#); [Dementieva et al., 2023](#))
 - In-context Learning ([Chen et al., 2024](#); [Zhang et al., 2024](#); [Pan et al., 2024](#))
 - Prompt-based text editing ([Liu et al., 2024](#); [Luo et al., 2023](#))

In this paper, we focus on a novel **decoding approach** to guide LLMs for TST using fixed prompts and therefore it does not require significant computational consumption and ensures stable outputs.

Neuron Analysis

- Multilinguality analysis (Kojima et al., 2024; Tan et al., 2024; Tang et al., 2024)
- Knowledge enhancement (Li et al., 2023)
- Sentiment analysis (Tigges et al., 2023)

Motivated by the promising outcomes of neuron analysis in enhancing multilingual capabilities of LLMs, this paper posits the presence of **style-specific neurons**, identifies them, and integrates neuron activation and deactivation seamlessly into the decoding process.

① Introduction

② Related Work

③ Method

④ Experiments

⑤ Results

⑥ Analysis

⑦ Conclusion

Overview

Identifying Style-Specific Neurons



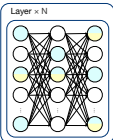
Positive

Positive Texts

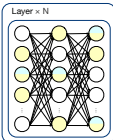


Negative

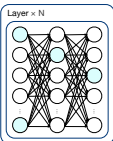
Negative Texts



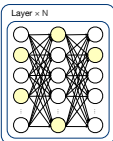
Positive activation neurons



Negative activation neurons



Positive selected neurons



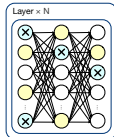
Negative selected neurons

○ Positive neuron ● Negative neuron ● Overlap neuron

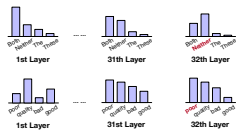
Deactivating Source Style Neurons

Instruction: Please transfer the following positive style sentence into a negative style sentence and maintain the meaning of the sentence.

Input: Both dishes were prepared with **quality** veggies.

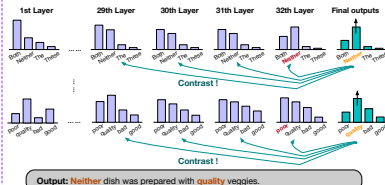


⊗ Deactivated neuron



Output: Neither dishes were prepared with **poor** veggies.

Decoding by Contrasting Style Layer



Output: Neither dish was prepared with **quality** veggies.

- Identify Style-Specific Neurons
 - Neuron Selection (remove overlap)
- Deactivating Source Style Neurons
 - Deactivate neurons from source side or target side
- Contrastive Decoding for TST
 - Adapt Dola ([Chuang et al., 2024](#)) to TST to mitigate the fluency issues observed during neuron deactivation

Identify Style-Specific Neurons

- Neurons in LLMs

$$a^{(j)} = \text{act_fn}(W^{(j)} a^{(j-1)} + b^{(j)})$$

- where $W^{(j)}$ and $b^{(j)}$ are the weights and biases of layer j , while $a^{(j-1)}$ is the activation values of the previous layer and $\text{act_fn}(\cdot)$ denotes the activation function (e.g., GLU used in LLaMA).
- The i^{th} neuron of the layer is considered to be active when its activation value $a_i^{(j)} > 0$.

Neuron Selection

• Overlap Analysis

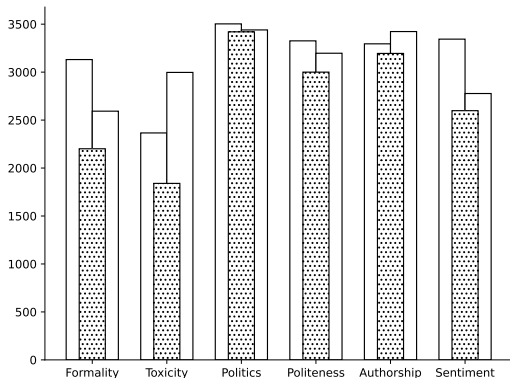


Fig. Overlap statistics of neuron using (Tang et al., 2024)

- Higher overlap among style-specific neurons when using traditional neuron selection method (e.g., 95% overlap in politics style).
- The overlap negatively impacts the performance of TST (details in ablation study).
- We **remove the intersection** between source style neuron and target style neuron.

Deactivating Source Style Neurons

- Deactivate the neurons by setting the activation values to zero during the model's forward pass.
- But, which side should we deactivate? Source or Target?

Style Accuracy					
Source	Target	Formality		Politeness	
		informal	formal	impolite	polite
✗	✗	80.00	11.20	79.50	14.80
✓	✗	80.53	13.63	80.06	19.37
✗	✓	76.25	8.51	65.50	9.27
✓	✓	78.42	9.27	73.48	10.36

Fluency					
Source	Target	Formality		Politeness	
		informal	formal	impolite	polite
✗	✗	92.53	87.69	105.35	92.34
✓	✗	104.17	96.83	127.26	105.12
✗	✓	113.14	106.23	136.10	112.51
✓	✓	108.22	100.79	131.22	108.64

Tab. Experiments for deactivating neurons on formality and politeness benchmarks.

- Deactivating the source-style neurons while keeping the targetstyle neurons active improves the accuracy of generating the target style.
- Fluency decreases whenever neurons are deactivated, whether they are source-style or target-style neurons.

Contrastive Decoding for TST

- Overview of Dola (Chuang et al., 2024)

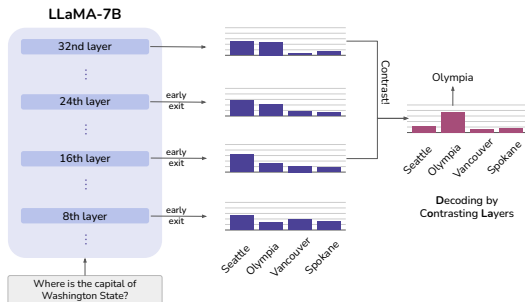


Fig. Figure directly taken from Chuang et al., 2024

- When predicting factual information, LLaMA tends to change the predictions in the higher layers. Otherwise, predictions usually have been decided by early layers.
- Three Steps:
 - Early exiting from all layers;
 - Pick a layer as “premature” layer, final layer as “mature” layer;
 - Subtract “premature” logits from “mature” logits in log domain.

Contrastive Decoding for TST

- Candidate layer selection

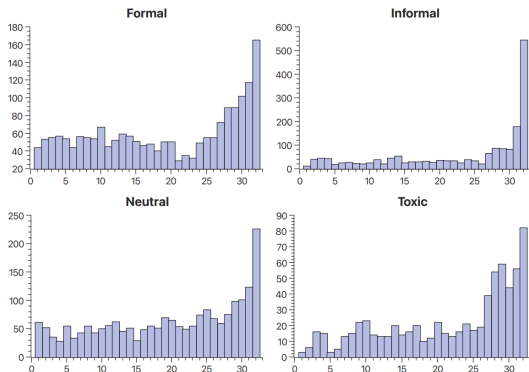


Fig. Statistics of the number of style-specific neurons in each layer

- The last few layers, particularly the final layer, contain significantly more style neurons compared to the earlier layers.
- We select the last few layers (4 in our experiments) as our candidate lay.

Contrastive Decoding for TST

- Next-token prediction
 - Select the layer with the maximum JSD distance from the candidate layers as our premature layer M and adjust their probability distribution.

① Introduction

② Related Work

③ Method

④ Experiments

⑤ Results

⑥ Analysis

⑦ Conclusion

Datasets

- We evaluate our approach on 6 TST benchmarks

Benchmark	Dataset	Tasks	Size		
			train	vald	test
Politeness	Politnss (Madaan et al., 2020)	impolite ↔ polite	100k	2000	2000
Toxicity	ParaDetox (Logacheva et al., 2022)	toxic ↔ neutral	18k	2000	2000
Formality	GYAFC (Rao and Tetreault, 2018)	informal ↔ formal	52k	500	500
Authorship	Shakespeare (Xu et al., 2012)	shakespeare ↔ modern	27k	500	500
Politics	Political (Voigt et al., 2018)	democratic ↔ republican	100k	1000	1000
Sentiment	Yelp (Shen et al., 2017)	positive ↔ negative	100k	1000	1000

Baselines

- We compare our approach with the following baselines:
 - **LLaMA-3:** We use LLaMA-3 ([Meta 2024](#)) without additional fine-tuning as the vanilla baseline system.
 - **APE:** Using activation probability entropy to identify the style specific neurons ([Tang et al., 2024](#)).
 - **AVF:** Using activation value frequency and set a threshold to identify the style neurons ([Tan et al., 2024](#)).
 - **PNMA:** Finding neurons that activate on the source style sentences but do not activate on target style sentences ([Kojima et al., 2024](#)).

Evaluation Metrics

- We evaluate our approach using three common metrics:
 - **Style Accuracy:** Accuracy of labels predicted as correct by a style classifier.
 - **Content Preservation:** Cosine similarity between the embeddings of the original text and the text generated by the model, using LaBSE (Feng et al., 2022) to obtain sentence embeddings as our primary metric.
 - **Fluency:** Perplexity of the generated sentences using GPT-2 (Radford et al., 2019).
- Classifiers used to evaluate the accuracy of style transfer

Benchmark	Source
Politeness	https://huggingface.co/Genius1237/xlm-roberta-large-tydip
Toxicity	https://huggingface.co/s-nlp/roberta_toxicity_classifier
Formality	https://huggingface.co/s-nlp/xlmr_formality_classifier
Authorship	https://huggingface.co/notaphoenix/shakespeare_classifier_model
Politics	https://huggingface.co/m-newhauser/distilbert-political-tweets
Sentiment	https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english

① Introduction

② Related Work

③ Method

④ Experiments

⑤ Results

⑥ Analysis

⑦ Conclusion

① Introduction

② Related Work

③ Method

④ Experiments

⑤ Results

⑥ Analysis

⑦ Conclusion

Ablation Study

- Neuron overlapping ablation analysis

	Style	without	with
Formality	informal→formal	74.00	79.40
	formal→informal	12.20	13.63
Toxicity	toxic→neutral	47.57	49.78
	neutral→toxic	28.44	29.82
Politics	democratic→republican	40.90	37.51
	republican→democratic	44.80	49.70
Politeness	impolite→polite	77.10	80.10
	polite→impolite	18.20	21.73
Authorship	shakespeare→modern	55.80	63.00
	modern→shakespeare	44.60	45.42
Sentiment	positive→negative	78.90	79.75
	negative→positive	48.00	51.70

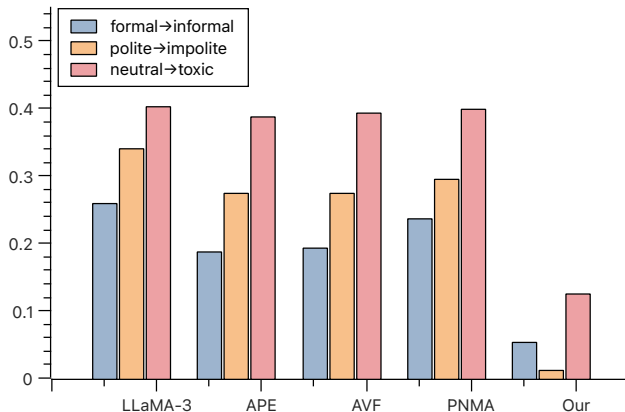
Table 3: **Ablation study:** Style transfer accuracy on removing overlap between source- and target-side style neurons in six benchmarks. “with” indicates the removal of overlap.

- Deactivation and contrastive decoding ablation analysis

	Deactivate	Contrastive	Toxicity		Authorship	
			toxic	neutral	shakespeare	modern
#1	✗	✗	47.67	29.04	63.80	43.80
#2	✓	✗	52.63	31.07	68.39	44.71
#3	✗	✓	46.82	28.31	63.23	43.16
#4	✓	✓	55.36	31.98	73.40	45.14

Table 4: **Ablation study:** Style transfer accuracy for neuron deactivation and contrastive decoding on the toxicity and authorship tasks. “✓” means the inclusion of the neuron deactivation or contrastive decoding steps, while “✗” means they are turned off. #1 indicates the results from baseline LLaMA-3 model, which do not use the deactivation nor the contrastive steps.

Copy Problem



- Neuron-based approaches (APE, AVF, and PNMA) partially mitigate this issue by controlling neuron activation;
- our approach achieves a reduced copy rate by deactivating source-side neurons and employing a novel decoding strategy.

Further Analysis

- Different model
 - Our method consistently demonstrates effectiveness across diverse model sizes, including larger models like 70B.
- Layer selection strategies
 - Selecting the last few layers proves optimal compared to earlier layers.
- Content preservation metrics
 - Different strategies for preserving meaning yield similar outcomes, highlighting the importance of exploring innovative approaches in future research.
- Decoding strategies
 - Contrastive decoding exhibits significant advantages over traditional decoding methods in the TST task, motivating our adoption of CD strategy.

① Introduction

② Related Work

③ Method

④ Experiments

⑤ Results

⑥ Analysis

⑦ Conclusion

Conclusion

- (i) To the best of our knowledge, this is the first work on using style-specific neurons to steer LLMs in performing text style transfer tasks.
- (ii) We emphasize the significance of eliminating overlap between neurons activated by source and target styles, a methodological innovation with potential applications beyond style transfer.
- (iii) We introduce an enhanced contrastive decoding method inspired by Dola. Our approach not only increases the production of words in the target style but also ensures the fluency of the generated sentences, addressing issues related to direct copying of input text in TST.

Thank You!

Email: wen.lai@tum.de

Homepage: <https://wenlai-lavine.github.io>

Address: Bildungscampus 3, 74076 Heilbronn, Germany



Paper



Code