

Controllable Text Generation for Style and Content

Michele Papucci * <>

* ItaliaNLP LAB @ CNR-ILC Oniversità di Pisa michele.papucci@phd.unipi.it



Introduction, Motivation, and Goal

Nowadays, Large Language Models (LLMs) have become the standard state-of-the-art technique to solve most Natural Language Processing (NLP) tasks, especially in text generation tasks. While extremely good at solving those, LLMs still suffer from major problems like following instructions or constraints, generating unfaithful or untrue texts (hallucinations), and their decisionmaking process being unknown. My project aims at solving these problems, and it's focused on using Controllable Text Generation (CTG) and Interpretability techniques to make the model adhere to specific syntactic or style-related constraints and to detect and prevent the generation of hallucinations.

Here are presented two works that use two different CTG techniques (a Decoding technique and a Reinforcement Learning technique) to control the model's output. In particular, the focus is on constraining the model to produce text that aligns with desired linguistic patterns.

Controllable Text Generation Experiments

Shifting Model Writing Style to Fool Detectors

We defined a CTG pipeline that by fine-tuning language models with Direct Preference Optimization (DPO), shifts the style of MGTs to resemble Human-written Texts (HWTs), making them harder to detect.



We present two techniques to build a preference dataset for a DPO fine-tuning:

- dpo: We create a parallel dataset of (HWT, MGT), and label the **HWT as the preferred option**.
- dpo-ling: We train a Linear SVM to solve the MGT detection task on the text linguistic profiling obtained with ProfilingUD. Then, we take the top-k features that have the highest absolute coefficient for the SVM classifica-

Generating Multi-level Text Simplification

LLMs can be used as a tool for generating datasets for lowresource languages. We propose a pipeline for the creation and evaluation of a Sentence Simplification resource for Italian:

- 1. We identified the best-performing Italian LLM for Sentence Simplification (LlaMantino-2);
- 2. We sampled simplifications from the model using original sentences from two different domains (Wikipedia and PaWaC) using the **Diverse Beam Search Decoding** technique, ensuring we obtained 10 different simplifications for each sentence;
- 3. We evaluated the resulting sentence pairs in terms of their linguistic feature diversity and variation in readability levels.

	Wikipedia	Pawac
	Pillai's Trace	Pillai's Trace
Driginal vs Least Simplified	.12	.16
Driginal vs Randomly-Selected	.18	.19
Driginal vs Most Simplified	.44	.46

We found that the selected LLM was capable of generating multiple simplifications, at different readability levels (evaluated with Read-IT), for each original sentence.

By looking at the linguistic profiling of the generated simplifications, we found that the model **used simplification**

We find various linguistic patterns shared between the two domains that are correlated with the text complexity:

• Sequence length is the most important proxy of readability;

- The use of **subordination** is highly relevant;
- The presence of **low-frequency words** impacts the sentence readability heavily.

Spearman Rank with Read-IT

	Wikipedia	PaWaC
tokens_per_sent -	0.51	0.42
average_class -	-0.06	
highest_class -	0.16	0.20
in_AD -		0.05
in_AD_types -	-0.16	-0.15
in_AU -		0.09
in_AU_types -		0.09
in EQ types		-0.08
in_ru_types -		-0.07
ava verh edaes	0.07	0.11
verb edaes dist 1 -	0.11	0.11
verb edges dist 2 -		0.06
verb_edges_dist 4 -	0.05	
verb_edges_dist_5 -	0.07	0.12
verb_edges_dist_6 -	0.07	0.07
verbal_head_per_sent -	0.39	0.28
verbal_root_perc	-0.08	
dep_dist_acl -	0.13	0.16
dep_dist_acl:relcl -	0.16	0.15
dep_dist_advcl -	0.11	0.09
dep_dist_advmod -	0.09	0.11
dep_dist_appos -	0.11	0.15
dop_dist_aux -	-0.06	-0.07
dop_dict_cocc	0.07	-0.05
den dist co	0.07	0.10
den dist comp	0.10	0.10
den dist compound	0.10	0.08
dep dist coni -	0.07	0.12
dep dist det -	-0.10	-0.11
dep_dist_det:poss -		0.08
dep_dist_expl -	0.09	0.11
dep_dist_expl:impers -	-0.05	
dep_dist_fixed -	0.08	0.06
dep_dist_flat:foreign -	0.05	
dep_dist_iobj -	0.07	
dep_dist_mark -	0.09	0.11
dep_dist_nmod -		0.06
dep_dist_nsubj -	-0.10	-0.05
dep_dist_nummod -		0.07
aep_aist_obl-	0.06	0.02
den dist parataxia	0.00	0.08
den dist nunct	0.09	0.13
dep dist xcomp -	0.07	0.06
avg links len -	0.30	0.30
avg_prepositional_chain_len -	0.14	0.16
avg_token_per_clause -	0.05	0.17
max_depth -	0.38	0.38
max_links_len -	0.32	0.31
n_prepositional_chains -	0.31	0.32
prep_dist_2 -	0.11	0.10
prep_dist_3 -		0.08
prep_dist_4 -		0.06
upos_dist_ADP -	0.07	0.10
upos_dist_ADV -	0.12	0.10
upos_dist_AUX -	-0.16	-0.16
upos_aist_CCONJ -	0.05	0.11
	-0.09	-0.11
upos_dist_NUUN -	-0.06	
upos_uisi_ivom -	0.15	0.18
upos dist PROPN -	-0.06	0120
upos dist PUNCT -	0.10	0.12
upos dist SCONI -	0.09	0.14
upos_dist_VERB -		-0.05
avg_subordinate_chain_len -	0.30	0.26
principal_proposition_dist -	-0.39	-0.24
subordinate_dist_1 -	0.17	0.18
subordinate_dist_2 -	0.11	0.09
subordinate_dist_3 -	0.10	0.06
subordinate_dist_4 -	0.05	
subordinate_post -	0.22	0.20
subordinate_pre -	0.09	0.10
suburumate_proposition_dist -	0.39	0.32
aux_10111_uist_1111 - aux form dict Part	-0.05	
aux mood dist Cnd	0.05	
aux mood dist Ind -	-0.05	
aux_mood_dist_Sub -		0.05
aux_tense_dist Imp -	0.06	0.04
aux_tense_dist_Pres -	-0.05	
verbs_form_dist_Fin -	0.05	0.06
verbs_form_dist_Ger -	0.13	0.06
verbs_form_dist_Inf -		0.06
verbs_form_dist_Part -		0.07
verbs_mood_dist_Ind -	0.17	0.13
verbs_mood_dist_Sub -	0.05	
verbs_num_pers_dist_Plur+3 -	0.08	0.07
verbs_num_pers_dist_Sing+1 -	0.10	0.05
verbs_num_pers_dist_Sing+3 -	0.13	0.12
verbs_tense_dist_Imp -		0.07
verbs_tense_dist_Past -	0.07	0.07
nhi nost	0.11	0.07
ohi nre -	0.07	0.06
subi nost -	0.05	0.10
)		

- 0.75

- 0.50

- 0.25

- 0.00

tion. For each of these features, we select pairs of (HWT, MGT) that **maximize that feature difference**, and tag the HWT as the preferred. The objective is to take the top-most representative couples of sentences for each of the most important features for the SVM detection.

	Performance Drop (F1) on Llama - XSUM				
	Mage	Radar	LLM-DetectAlve	Binoculars	
dpo	0.36	0.15	0.19	0.66	
dpo-ling	0.29	0.36	0.18	0.61	

What we found is that the *dpo* approach was the one that maximized the *performance drop* of the four state-of-the-art detectors we tested.

Instead, *dpo-ling* was the technique that performed better in the objective of aligning feature-specific distribution of the generated texts to the human writing style.

Human raters' performances were *unaffected* by the DPO fine-tuning, remaining at around 50% accuracy, the same as blindly guessing.



strategies similar to those used in manually crafted simplification resources. In particular, the two methodologies shared:

• a lower number of tokens;

- fewer pronouns, adverbs, and punctuation marks, a higher proportion of *determiners*;
- a shallower syntactic trees, and shorter dependency links.

With this newly created resource, we aim to train models that can be used to produce **simplification aimed at a spe**cific readability level of the user.



-0.25

-0.50

-0.75

Figure 1: The distribution of selected linguistic features on Llama XSUM generation.

Figure 2: For both Wikipedia on the left and PaWaC on the right, the Kernel Density Estimate for the READ-IT.

1.2

Figure 3: Correlation between linguistic feature differences (original vs. simplified) and READ-IT scores. White cells indicate non-statistically significant Spearman rank correlations (p-value < 0.05).