



HAL
open science

Sur les limites de l'identification par l'humain de textes générés automatiquement

Nadège Alavoine, Maximin Coavoux, Emmanuelle Esperança-Rodier, Romane Gallienne, Carlos-Emiliano González-Gallardo, Jérôme Goulian, José G. Moreno, Aurélie Névéol, Didier Schwab, Vincent Segonne, et al.

► To cite this version:

Nadège Alavoine, Maximin Coavoux, Emmanuelle Esperança-Rodier, Romane Gallienne, Carlos-Emiliano González-Gallardo, et al.. Sur les limites de l'identification par l'humain de textes générés automatiquement. 35èmes Journées d'Études sur la Parole (JEP 2024) 31ème Conférence sur le Traitement Automatique des Langues Naturelles (TALN 2024) 26ème Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues (RECITAL 2024), Jul 2024, Toulouse, France. pp.18-19. hal-04623002

HAL Id: hal-04623002

<https://inria.hal.science/hal-04623002v1>

Submitted on 28 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Sur les limites de l'identification par l'humain de textes générés automatiquement

Nadège Alavoine¹, Maximin Coavoux², Emmanuelle Esperança-Rodier²,
Romane Gallienne³, Carlos-Emiliano González-Gallardo⁴, Jérôme Goulian²,
Jose G. Moreno⁵, Aurélie Névéol⁶, Didier Schwab²,
Vincent Segonne⁷ and Johanna Simoens⁸

(1) Université Paris-Saclay, LISN, Campus Universitaire bâtiment 507, Rue du Belvédère, 91400 Orsay, France

(2) Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, 38000 Grenoble, France

(3) Université Sorbonne Nouvelle, Lattice, CNRS, ENS-PSL, 1 rue Maurice Arnoux, 92120 Montrouge, France

(4) La Rochelle Université, L3i, 17000 La Rochelle, France

(5) University of Toulouse, IRIT, 31000 Toulouse, France

(6) LISN, Université Paris-Saclay, CNRS, 91403 Orsay, France

(7) Université Bretagne Sud, UMR CNRS 6074, IRISA, F-56000 Vannes, France

(8) Everteam, Bagneux, France

nadège.alavoine@universite-paris-saclay.fr,

{first.last}@univ-grenoble-alpes.fr

romane.gallienne@cnrs.fr, carlos.gonzalez_gallardo@univ-lr.fr,

jose.moreno@irit.fr, aurelie.neveol@lisn.upsaclay.fr

vincent.segonne@univ-ubs.fr, johanna.simoens@gmail.com

RÉSUMÉ

La génération de textes neuronaux fait l'objet d'une grande attention avec la publication de nouveaux outils tels que ChatGPT. La principale raison en est que la qualité du texte généré automatiquement peut être attribuée à un-e rédacteurice humain-e même quand l'évaluation est faite par un humain. Dans cet article, nous proposons un nouveau corpus en français et en anglais pour la tâche d'identification de textes générés automatiquement et nous menons une étude sur la façon dont les humains perçoivent ce texte. Nos résultats montrent, comme les travaux antérieurs à l'ère de ChatGPT, que les textes générés par des outils tels que ChatGPT partagent certaines caractéristiques communes mais qu'ils ne sont pas clairement identifiables, ce qui génère des perceptions différentes de ces textes par l'humain. Ceci est le résumé de l'article "Limitations of Human Identification of Automatically Generated Text" publié à LREC-COLING-2024 (Alavoine *et al.*, 2024).

ABSTRACT

Here the title in English.

Neural text generation is receiving broad attention with the publication of new tools such as ChatGPT. The main reason for that is that the achieved quality of the generated text may be attributed to a human writer by the naked eye of a human evaluator. In this paper, we propose a new corpus in French and English for the task of recognising automatically generated texts and we conduct a study of how humans perceive the text. Our results show, as previous work before the ChatGPT era, that the generated texts by tools such as ChatGPT share some common characteristics but they are not clearly identifiable which generates different perceptions of these texts.

MOTS-CLÉS : identification humaine, génération de texte avec des modèles neuronaux, ChatGPT.

Références

ALAVOINE N., COAVOUX M., ESPERANÇA-RODIER E., GALLIENNE R., GALLARDO C. G., GOULIAN J., MORENO J. G., NEVEOL A., SCHWAB D., SEGONNE V. *et al.* (2024). Limitations of human identification of automatically generated text. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, p. 10511–10516.