

Pantagruel : des encodeurs auto-supervisés unifiés pour le texte et la parole

Phuong-Hang Le^{1,9} Valentin Pelloin² Arnault Chatelain⁴ Maryem Bouziane³
Mohammed Ghennai¹ Qianwen Guan⁵ Kirill Milintsevich²
Salima Mdhaffar³ Aidan Mannion¹ Nils Defauw⁶ Shuyue Gu⁵
Alexandre Audibert¹ Marco Dinarelli¹ Yannick Estève³ Lorraine Goeriot¹
Steffen Lalande² Nicolas Hervé² Maximin Coavoux¹ François Portet¹
Étienne Ollion⁴ Marie Candito⁵ Maxime Peyrard¹ Solange Rossato¹
Benjamin Lecouteux¹ Aurélie Nardy⁷ Gilles Sérasset¹ Vincent Segonne⁸
Solène Evain¹⁰ Diandra Fabre¹ Didier Schwab¹

¹ Univ. Grenoble Alpes, CNRS, Grenoble INP, LIG, 38000 Grenoble, France

² INA (Institut National de l’Audiovisuel), 4 Avenue de l’Europe, 94366 Bry-sur-Marne, France

³ Avignon Université, LIA, France

⁴ CREST (École Polytechnique, ENSAE, CNRS), 5 avenue Le Chatelier, 91120 Palaiseau, France

⁵ LLF (Université Paris Cité and CNRS), UFRL Olympe de Gouges,

13 place Paul Ricoeur, 75013 Paris, France

⁶ Univ. Grenoble Alpes, EFELIA-MIAI, IUT2 Grenoble, LIG, 38000 Grenoble, France

⁷ Univ. Grenoble Alpes, Lidilem, 38000 Grenoble, France

⁸ Université Bretagne Sud, CNRS, IRISA, France

⁹ Saclay AI, France

¹⁰ IRIT, Université de Toulouse, CNRS, Toulouse INP, UT3, Toulouse, France

didier.schwab@univ-grenoble-alpes.fr

RÉSUMÉ

Nous publions les modèles Pantagruel, une nouvelle famille d’encodeurs autosupervisés pour le texte et la parole en français. Plutôt que de prédire des cibles adaptées à chaque modalité, comme des tokens textuels ou des unités de parole, Pantagruel apprend des représentations cibles contextualisées dans l’espace des caractéristiques, ce qui permet à des encodeurs spécifiques à chaque modalité de mieux capturer les régularités linguistiques et acoustiques. Des modèles distincts sont préentraînés sur de vastes corpus français, notamment Wikipedia, OSCAR et CroissantLLM pour le texte, ainsi que MultilingualLibriSpeech, LeBenchmark et INA-100k pour la parole. INA-100k est un nouveau corpus de 100 000 heures d’audio en français, issu des archives de l’Institut national de l’audiovisuel (INA), dépôt national des émissions de radio et de télévision françaises, fournissant des données audio très diversifiées. Nous évaluons Pantagruel sur un large éventail de tâches aval couvrant les deux modalités, y compris celles des principaux benchmarks français tels que FLUE ou LeBenchmark. Sur l’ensemble de ces tâches, les modèles Pantagruel obtiennent des performances compétitives, voire supérieures, par rapport à de solides bases de référence françaises comme CamemBERT, FlauBERT et LeBenchmark 2.0, tout en conservant une architecture commune capable de traiter de manière transparente des entrées de parole ou de texte. Ces résultats confirment l’efficacité d’objectifs auto-supervisés dans l’espace des caractéristiques pour l’apprentissage de représentations en français et mettent en évidence Pantagruel comme une base robuste pour la compréhension multimodale

parole-texte.

ABSTRACT

Pantagruel : Unified Self-Supervised Encoders for French Text and Speech

We release Pantagruel models, a new family of self-supervised encoder models for French text and speech. Instead of predicting modality-tailored targets such as textual tokens or speech units, Pantagruel learns contextualized target representations in the feature space, allowing modality-specific encoders to capture linguistic and acoustic regularities more effectively. Separate models are pre-trained on large-scale French corpora, including Wikipedia, OSCAR and CroissantLLM for text, together with MultilingualLibriSpeech, LeBenchmark, and INA-100k for speech. INA-100k is a newly introduced 100000-hours corpus of French audio derived from the archives of the Institut National de l'Audiovisuel (INA), the national repository of French radio and television broadcasts, providing highly diverse audio data. We evaluate Pantagruel across a broad range of downstream tasks spanning both modalities, including those from the standard French benchmarks such as FLUE or LeBenchmark. Across these tasks, Pantagruel models show competitive or superior performance compared to strong French baselines such as CamemBERT, FlauBERT, and LeBenchmark 2.0, while maintaining a shared architecture that can seamlessly handle either speech or text inputs. These results confirm the effectiveness of feature-space self-supervised objectives for French representation learning and highlight Pantagruel as a robust foundation for multimodal speech–text understanding.

MOTS-CLÉS : apprentissage auto-supervisé, JEPa, data2vec, modèles de langue français, encodeurs parole et texte, apprentissage de représentations multimodales, architecture prédictive à embeddings joints, modélisation prédictive.

KEYWORDS: self-supervised learning, JEPa, data2vec, French language models, speech and text encoders, multimodal representation learning, joint-embedding predictive architecture, predictive modeling.

ARTICLE ACCEPTÉ À : LREC 2026 : The Fifteenth biennial Language Resources and Evaluation Conference, Palma, Mallorca, Spain, May 11-16, 2026. .

URL : <https://lrec.elra.info/lrec2026-main-799>

Remerciements

Ces travaux de recherche ont été partiellement financés par l'Agence nationale de la recherche (ANR), projet « PANTAGRUEL », ANR-23-IAS1-0001. Ils ont également bénéficié du soutien du projet CREMA (Coreference REsolution into MACHine translation), financé par l'ANR, contrat n° ANR-21-CE23-0021-01. Ils ont par ailleurs reçu un financement public géré par l'ANR dans le cadre de la stratégie France 2030, référence ANR-23-IACL-0006. Ces travaux ont également été soutenus par l'ANR à travers la chaire MIAI « IA et Langage » (ANR-19-P3IA-0003), la chaire MIAI « Socialisation et Langage à l'école », et la chaire AugmentIA sous le mécénat d'Artelia et de Grenoble INP (ANR-23-IACL-0006).

Ce travail a été réalisé en utilisant les ressources HPC de GENCI à IDRIS et CINES dans le cadre des allocations 2022-A0131013801, 2023-A0151013801, 2024-A0171013801, 2024-A0161015074 et 2025-A0191013801 sur les supercalculateurs Jean Zay et AdastrA.