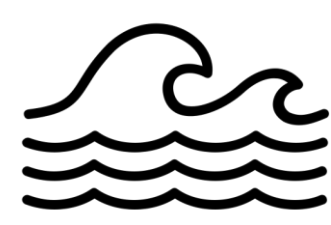# 3rd Edition
# Data Research meetup by MagIC

## WAVe:
## Word-Aligned Verification of Synthetic Speech for Automatic Speech Recognition

**Yuriy Perezhohin ( yperezhohin@novaims.unl.pt )**

NOVA Information Management School (NOVA IMS), Universidade Nova de Lisboa, Campus de Campolide, 1070-312 Lisboa, Portugal;

REMYND, Alameda Bonifácio Lázaro Lozano nº15, 1ºC, 2780-125, Oeiras, Portugal

## INTRODUCTION

Automatic Speech Recognition (ASR) for low-resource languages often relies on synthetic speech generated by pairing Large Language Model (LLM) transcripts with Text-to-Speech (TTS) generated audio. However, blindly incorporating synthetic data can introduce mispronunciations, word omissions, and prosodic anomalies that degrade ASR performance and increase training time [1], [2]. Previous filtering methods assessed audio-text similarity at the sentence level, which masks localized word-level errors, a synthetic utterance may appear semantically aligned while still containing critical defects [3].
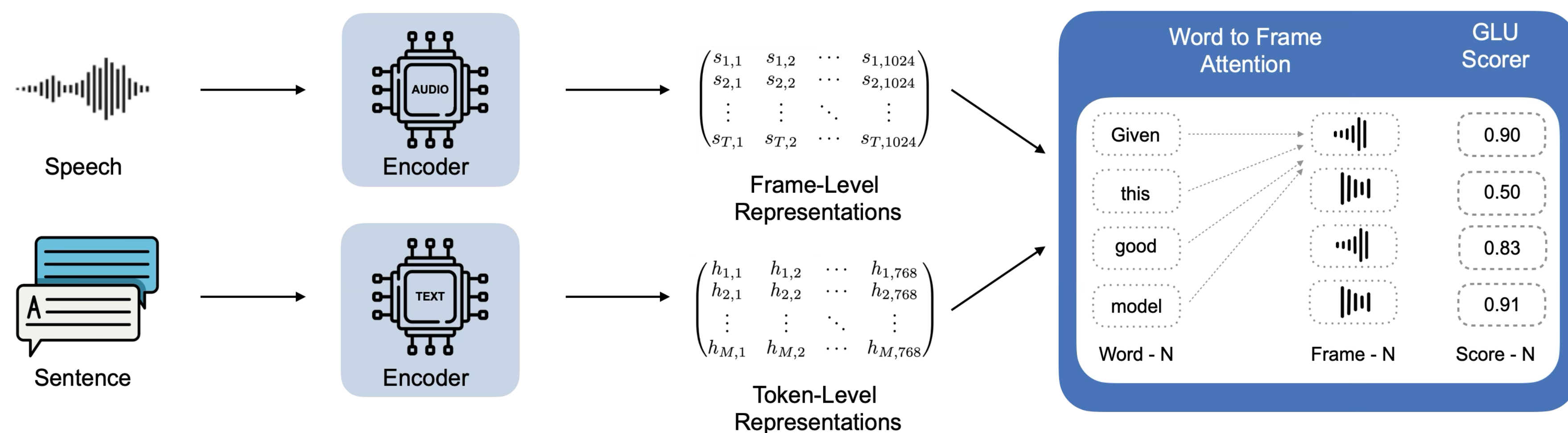
We present **WAVe** (Word-Aligned Verification), a multimodal embedding model that verifies each word against its corresponding audio frame through attention-based alignment. A Gated Linear Unit (GLU) scorer assigns confidence scores to each word, enabling fine-grained quality assessment.
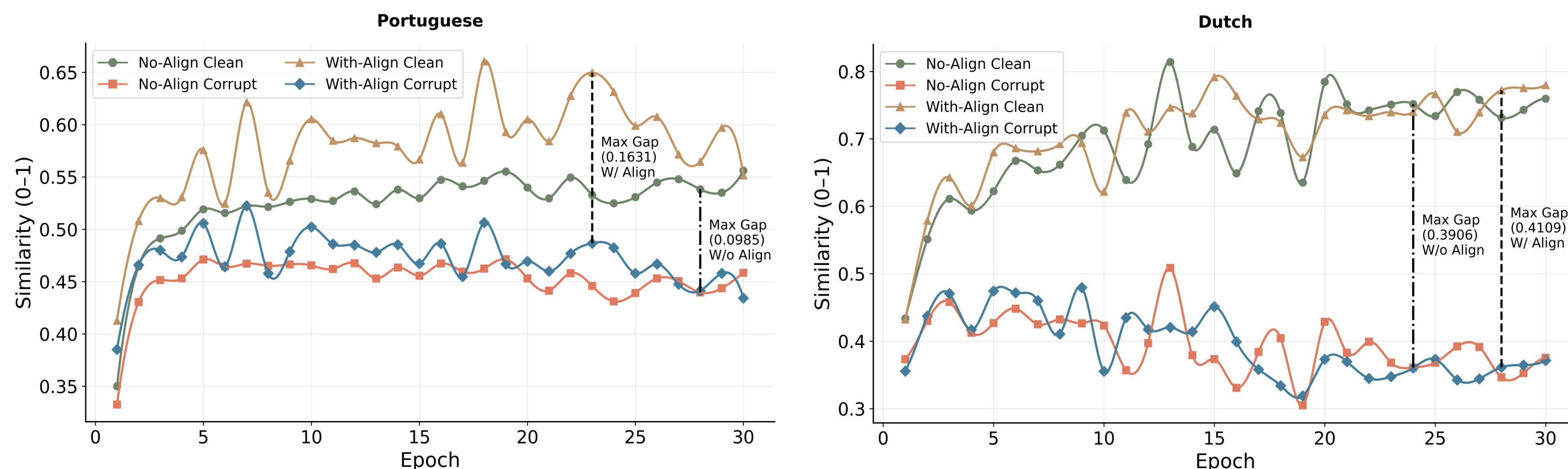
## METHODS AND MATERIALS

Our methodology follows a four-stage pipeline:
1. *WAVe Training* — We train WAVe via contrastive learning on Common Voice, contrasting clean audio-text pairs with corrupted versions to learn alignment.
2. *Synthetic Generation* — We then generate synthetic data using GPT-4o-mini transcripts and OpenAI's TTS (22k Portuguese, 35k Dutch samples).
3. *Quality Filtering* — The trained WAVe scores each pair ($q \in [0,1]$), partitioning them into subsets as shown in the table above.
4. *ASR Evaluation* — Finally, we fine-tune Whisper (Tiny, Small, Large-v3) on filtered data and evaluate on Common Voice and MLS benchmarks.

The WAVe architecture (right) illustrates the flow from encoders through cross-modal attention to word-level confidence scoring.



## RESULTS & DISCUSSION



### WAVe Training Performance

Word-level alignment supervision notably improves the model's ability to distinguish clean from corrupted pairs. For Portuguese, figure on the left, the clean-corrupt similarity gap increases from 9.85% (without alignment) to 16.31% (with alignment), a 6.5% absolute improvement achieved 5 epochs earlier. The Dutch language, figure on the right shows larger absolute margins (41.09%gap) due to its 60% larger training corpus.



### ASR Performance

WAVe demonstrates strong performance for Whisper Large-V3 on Portuguese across both evaluation benchmarks. On Common Voice (in-domain), high-quality filtering achieves 7.94% WER compared to 8.33% with unfiltered data, while requiring only 575 training steps, a 33% reduction. This efficiency stems from removing low-quality samples that introduce conflicting acoustic signals, allowing faster convergence. The benefits become even more pronounced for cross-domain generalization on MLS. WAVe-filtered CAPES dataset reduces WER from 13.54% to 6.89%, using 19% fewer steps and 30% less data. This substantial gain indicates that word-level filtering retains samples with robust acoustic representations that transfer well across domains, rather than overfitting to synthesis artifacts and hindering performance.

## CONCLUSION

Our approach achieves a **34% average reduction** in training steps while simultaneously decreasing WER, demonstrating that quality-based filtering delivers both computational efficiency and improved accuracy. Most critically, WAVe exhibits strong robustness: cross-domain evaluation reveals up to 49% improvement in generalization (13.54% → 6.89% MLS WER) using 30% less data. By retaining only high-quality synthetic samples, our method ensures models learn transferable acoustic representations, establishing quality-over-quantity as the key paradigm for efficient low-resource ASR augmentation.

## REFERENCES

- [1] Dhamyal et al. (2024). Using Voicebox-based Synthetic Speech for ASR Adaptation. *Proc. SynData4GenAI*.
- [2] Wang et al. (2025). From Tens of Hours to Tens of Thousands: Scaling Back-Translation for Speech Recognition. *arXiv*.
- [3] Manco et al. (2022). Learning music audio representations via weak language supervision. ICASSP 2022.

Funded by:

NOVA IMS MagIC · NOVA IMS · FCT Fundação para a Ciência e a Tecnologia · PRR · REPÚBLICA PORTUGUESA · Financiado pela União Europeia NextGenerationEU

NOVA Information Management School, 18 December 2025

**ONLINE VERSION**