

# **Accurate Detection of RNA Modifications on Nanopore Sequencing Data** at Signal-level with Machine Learning

Servi, L. <sup>1,2</sup>; Petrillo, E. <sup>1,2</sup>; Stegmayer G. <sup>3</sup>.

<sup>1</sup> Universidad de Buenos Aires, Facultad de Ciencias Exactas y Naturales, Departamento de Fisiología, Biología, Molecular, y Celular, Buenos Aires, Argentina

<sup>2</sup> CONICET-Universidad de Buenos Aires, Instituto de Fisiología, Biología Molecular y Neurociencias (IFIBYNE), C1428EHA, Buenos Aires, Argentina <sup>3</sup> Instituto de Investigación en Señales, Sistemas e Inteligencia Computacional, sinc(i), FICH-UNL/CONICET, Argentina

# INTRODUCTION

Accurate RNA modification identification is vital for understanding gene expression regulation and its implications in biotechnology. This study focuses on precise RNA modification detection using Oxford Nanopore Technology (ONT) sequencing data. RNA modifications play a central role in gene expression and post-transcriptional processes, underscoring their critical importance in diverse fields. Traditional methods often oversimplify by relying on reading errors as indirect indicators of RNA methylation, resulting in lower accuracy and a lack of specific modification details.

# **MATERIALS AND METHODS**

**ONT enables comprehensive analysis from its raw electrical signal** data. With Direct RNA Sequencing (DRS), numerous RNA features can be extracted. Each consumable flow cell contains 2048 pores actively incorporating nucleotide chains during approximately 48 hours of sequencing.

### PROPOSAL

Our approach focuses on the development of a Machine Learning (ML) model based on Natural Language Processing (NLP) models, harnessing the electrical current generated during sequencing and taking advantage of basecalling errors as fundamental features.





me-RIP sequencing

#### For training, we will employ:

- *me-RIPseq* data + Direct RNA Sequencing (DRS)
- Raw read signal from Nanopore



This device detects electric current changes from nucleic acids passing through the pore, measured in picoAmperes (pA) over time. A deep neural network, like Guppy, processes the signal to generate the nucleotide sequence in a step known as "basecalling."







#### CONCLUSIONS

Our proposal is a novel approach to the precise detection of RNA modifications with a ML model based on NLP and sequencing data. The accurate identification of RNA modifications will undoubtedly enhance our knowledge of gene regulation and open up new avenues for biotechnological applications.

Each covers pore approximately 5 nucleotides as they pass through, forming a signal from a k-mer instead of a single base.





- Pagès-Gallego, M., de Ridder, J. Comprehensive benchmark and architectural analysis of deep learning models for nanopore sequencing basecalling. Genome Biol 24, 71 (2023). https://doi.org/10.1186/s13059-023-02903-2
- Zhong ZD et al. Systematic comparison of tools used for m6A mapping from nanopore direct RNA sequencing. Nat Commun. 2023 Apr 5;14(1):1906. https://doi.org/10.1038/s41467-023-37596-5
- Matthew T Parker et al. Nanopore direct RNA sequencing maps the complexity of Arabidopsis mRNA processing and m6A modification eLife 9:e49658 (2020) https://doi.org/10.7554/eLife.49658
- https://github.com/nanoporetech/bonito
- https://distill.pub/2017/ctc/