# Evidence of Policy Violation by Reviewer wNfJ

Paper: #6307

## 1. Summary

This report presents evidence that the review provided by Reviewer wNfJ (Score: 2) was **fully generated using an LLM**, violating ICLR's review policy. Our analysis includes a control group: the same detection tools correctly identified all other reviewers (WqcR, rGrB, CKa4, F3q3) as "Fully Human," eliminating the possibility that the paper's technical jargon triggered a false positive. The summary of detection results is shown in Table 1.

Table 1: Summary of LLM detection result.

| Reviewer | Rating | Pangram Labs | GPTZero |
|---|---|---|---|
| wNfJ | 2 | Fully AI-generated | 100% AI generated |
| WqcR | 4 | | |
| rGrB | 6 | Fully human-written | 100% Human |
| CKa4 | 6 | | |
| F3q3 | 6 | | |

## 2. Evidence

### 2a. Pangram Labs

Pangram Labs (https://www.pangram.com/) is a leading enterprise-grade detection solution trusted by global organizations and academic institutions. Notably, Pangram provides a specialized detection module for ICLR reviews (https://iclr.pangram.com), specifically calibrated to handle the technical density and formal tone of ICLR submissions.

As shown in Figure 1, the tool successfully filtered our submission's reviews. Only **Reviewer wNfJ** (Score: 2) was flagged as **"Fully AI-generated"** (Red). The table reveals a perfect correlation between the AI flag and the low quality: the AI-generated review corresponds to the outlier score of 2.00, whereas the verified human reviews maintain a high average rating of 5.50.



Figure 1: Screen shot from Pangram Labs detection result for paper # 6307.

Official report link: https://iclr.pangram.com/reviews?submission_number=6307

## 2b. GPTZero

GPTZero (https://gptzero.me/) is widely regarded as the "gold standard" for AI detection, used by educators and institutions worldwide. It utilizes perplexity and burstiness metrics to distinguish between human and machine writing patterns.

Figure 2 identifies **Reviewer wNfJ as "100% AI-generated."** In contrast, Figures 3 and 4 confirm that reviews from other reviewers were correctly identified as "100% Human."
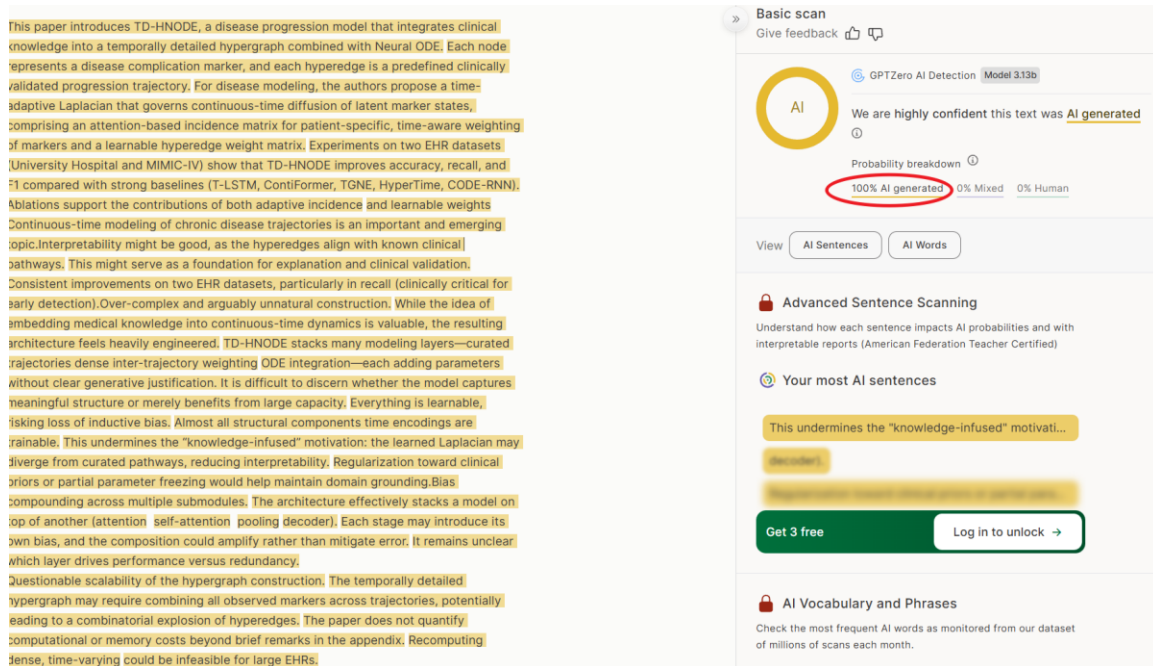


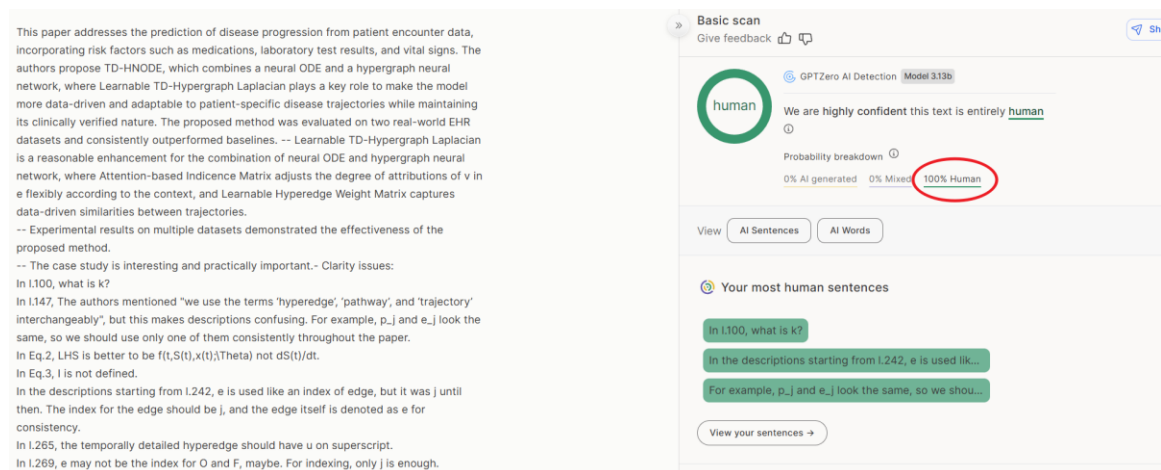Figure 2: GPTZero detection result (100% AI generated) for **Reviewer wNfJ.**



Figure 3: Control check: GPTZero detection result (100% Human) for Reviewer rGrB.
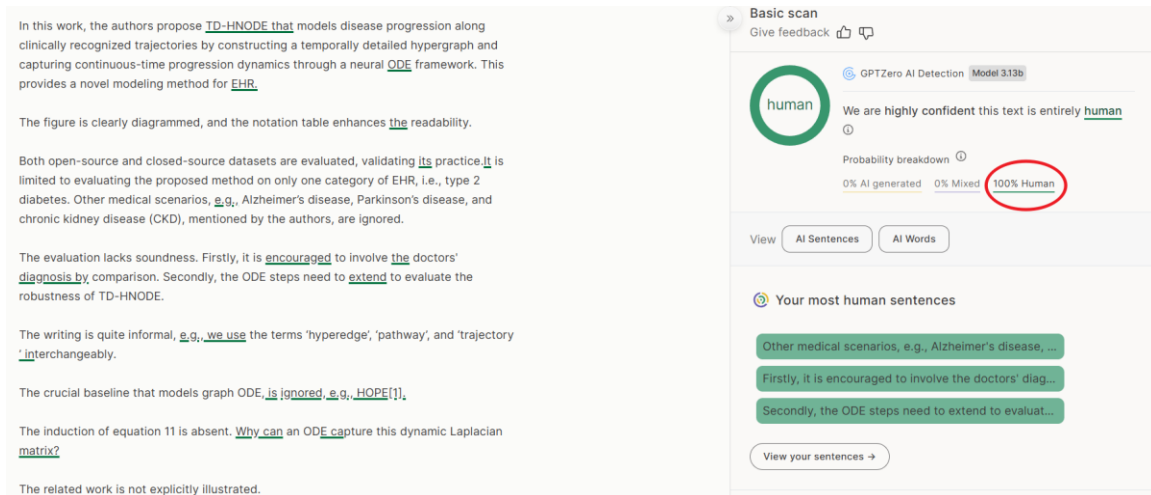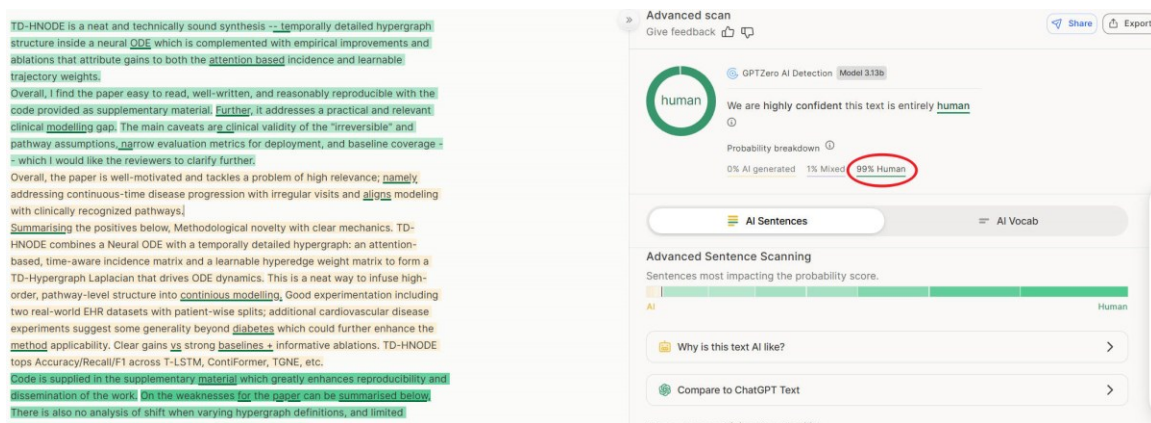
In this work, the authors propose TD-HNODE that models disease progression along clinically recognized trajectories by constructing a temporally detailed hypergraph and capturing continuous-time progression dynamics through a neural ODE framework. This provides a novel modeling method for EHR.

The figure is clearly diagrammed, and the notation table enhances the readability.

Both open-source and closed-source datasets are evaluated, validating its practice. It is limited to evaluating the proposed method on only one category of EHR, i.e., type 2 diabetes. Other medical scenarios, e.g., Alzheimer's disease, Parkinson's disease, and chronic kidney disease (CKD), mentioned by the authors, are ignored.

The evaluation lacks soundness. Firstly, it is encouraged to involve the doctors' diagnosis by comparison. Secondly, the ODE steps need to extend to evaluate the robustness of TD-HNODE.

The writing is quite informal, e.g., we use the terms 'hyperedge', 'pathway', and 'trajectory' interchangeably.

The crucial baseline that models graph ODE, is ignored, e.g., HOPE[1].

The induction of equation 11 is absent. Why can an ODE capture this dynamic Laplacian matrix?

The related work is not explicitly illustrated.

Figure 4: Control check: GPTZero detection result (100% Human) for Reviewer WqcR.



TD-HNODE is a neat and technically sound synthesis -: temporally detailed hypergraph structure inside a neural ODE which is complemented with empirical improvements and ablations that attribute gains to both the attention based incidence and learnable trajectory weights.
Overall, I find the paper easy to read, well-written, and reasonably reproducible with the code provided as supplementary material. Further, it addresses a practical and relevant clinical modelling gap. The main caveats are clinical validity of the "irreversible" and pathway assumptions, narrow evaluation metrics for deployment, and baseline coverage -- which I would like the reviewers to clarify further.
Overall, the paper is well-motivated and tackles a problem of high relevance; namely addressing continuous-time disease progression with irregular visits and aligns modeling with clinically recognized pathways.
Summarising the positives below, Methodological novelty with clear mechanics. TD-HNODE combines a Neural ODE with a temporally detailed hypergraph: an attention-based, time-aware incidence matrix and a learnable hyperedge weight matrix to form a TD-Hypergraph Laplacian that drives ODE dynamics. This is a neat way to infuse high-order, pathway-level structure into continuous modelling. Good experimentation including two real-world EHR datasets with patient-wise splits; additional cardiovascular disease experiments suggest some generality beyond diabetes which could further enhance the method applicability. Clear gains vs strong baselines + informative ablations. TD-HNODE tops Accuracy/Recall/F1 across T-LSTM, ContiFormer, TGNE, etc.
Code is supplied in the supplementary material which greatly enhances reproducibility and dissemination of the work. On the weaknesses for the paper can be summarised below,
There is also no analysis of shift when varying hypergraph definitions, and limited

Figure 5: Control check: GPTZero detection result (99% Human) for Reviewer CKa4.



[1] Learning the natural history of human disease with generative transformers. Nature 2025.

For linear Transformer, it can be considered as discreted ODE, for instance,

[2] TrajGPT: Irregular Time-Series Representation Learning of Health Trajectory. IEEE J-BHI 2025.

I just list few recent works. You could find more related works which should be included in the literature review.

Although this paper claims that it can generate interpretable trajectory, it has limited data analysis about the generated trajectories. It should include more case studies about how a patient's health progress over time. It could also include population-level analysis about the subphenotypes or combritidy.

It lacks interpretation and visualization of the learned token embedding (like HP, AF). We do not know whether model learns meaingful embedding or not in the graph-based model. Extending the weakness 2.

While the motivations mentions subpheontypes, did the author try to analyze it? whether it is connected the different progression speed in Fig.4.

This work focuses on T2D and selects features to analyze it. Did the author analyze the correlation or combritidy between T2D and other features? whether T2D will also contribute to other diseases (like heart failure)

While the background mentions medication and talks about "timely treatment", it does not have any thing about the interactions between disease progression and medications? we should see medication helps the disease recovery?
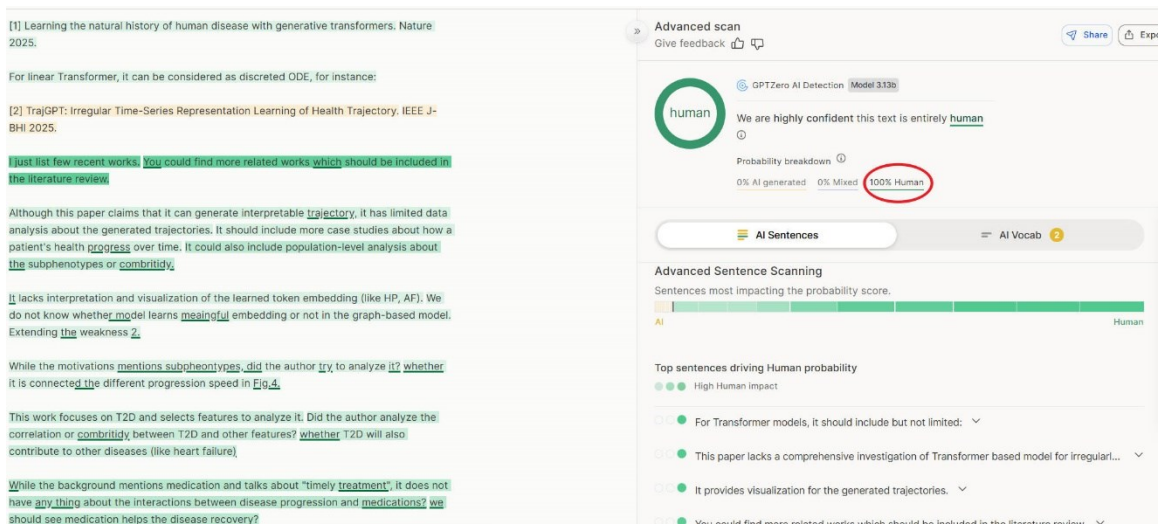
Figure 6: Control check: GPTZero detection result (100% Human) for Reviewer F3q3.