# Comparing Traditional and Graph-Enhanced Pipelines for Sepsis Prediction Using Omics Data

Joelle ASSY, Silvia TROTTET, Rayane ADAM, Yazid HOBLOS

Master 2 GENIOMHE-AI, Univ. Evry Paris-Saclay

**GRADUATE SCHOOL**
Informatique et Sciences du Numérique

## Introduction
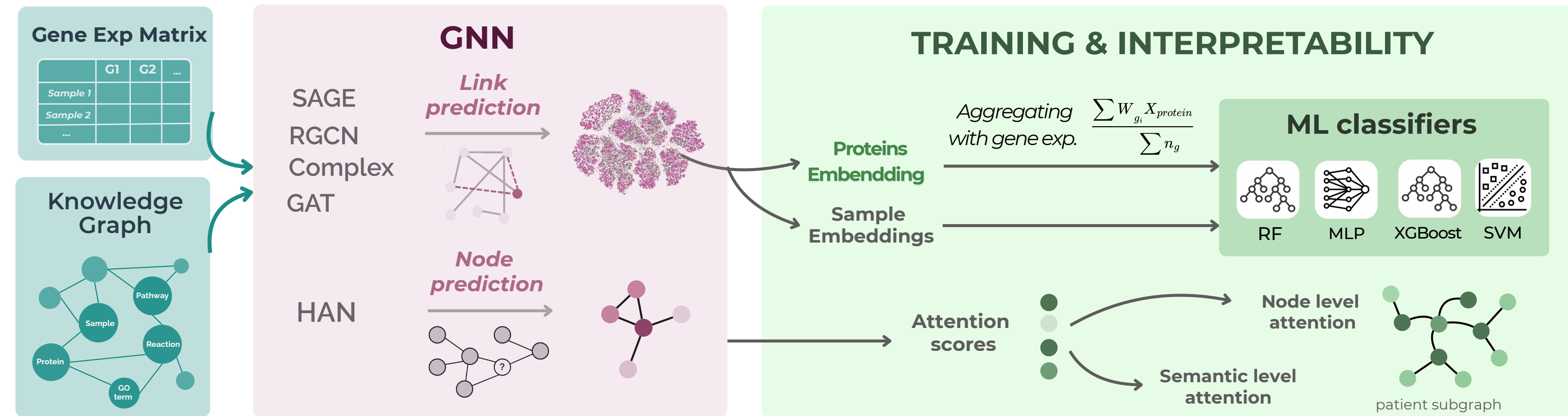
**Sepsis remains a leading cause of global mortality**

Sepsis is a life-threatening organ dysfunction resulting from a dysregulated host response to infection. Accurate early prediction is vital, yet traditional methods treat genes as independent features. This project investigates if outcomes are better predicted by encoding relational information among proteins, reactions, pathways, and GO terms. We compare a traditional transcriptomic pipeline against a graph-augmented approach that captures structural biological context.

**Datasets: (1) Gene Expression (2) Knowledge Graph**

**(1)** GSE54514 dataset with 163 whole blood samples from 26 survivors, 9 non-survivors, 18 controls collected daily for 5 days.

**(2)** Multi-relational network mapping Patient Samples, Proteins, GO Terms, Reactions, and Pathways, by integrating curated data from 8 sources: GO, Pathway Ontology, PPI, Entrez, BioGRID, UniProt, STRING, and Reactome.

## Overall Architecture



GNN: SAGE, RGCN, Complex, GAT — *Link prediction*; HAN — *Node prediction*

TRAINING & INTERPRETABILITY

*Aggregating with gene exp.* $\frac{\sum W_{g_i} X_{protein}}{\sum n_g}$

Proteins Embedding → ML classifiers (RF, MLP, XGBoost, SVM)
Sample Embeddings

Attention scores → Node level attention / Semantic level attention — patient subgraph

## Results

### ComplEx embeddings enhance sepsis prediction



**Graph protein embeddings are mostly better than gene expression across all metrics**

- ComplEx embeddings are top ranked
- ComplEx, GraphSAGE and GAT perform well for MLP
- SVM gene expression (93% acc) outperformed by ComplEx (95.3%)
- Recall of 100% and precision 77.6% are due to small dataset with imbalance
- RGCN protein embeddings are generally not great for prediction in different models

Accuracy over 10 seeds:
- ComplEx SVM top robust
- weighted RGCN nor reliable across models



### HAN patient-level and SHAP-ComplEx population-level explanations provide complementary insights



**HAN-based patient-specific explanation**
HAN-derived heterogeneous subgraphs highlight proteins with highest gradient-based influence on sepsis prediction, propagating influence through functional (GO) and pathway-level contexts to drive individual predictions.

**ComplEx-based consensus PPI drivers**
SHAP analysis across all ComplEx ML models yields a PPI network enriched in hub proteins, suggesting shared regulation.

## Discussion

- **Data scarcity:** dataset limited to 53 patients undermines training and lead to testing instability.
- **Class imbalance:** Only 22% healthy samples introduces bias toward septic phenotypes.
- **Temporal nature:** 163 samples spanning multiple time points require specialized modeling.
- **Unoptimized knowledge graph:** KG construction used no filtering, resulting in dense PPIs, pathway-dominated topology, and high reaction embeddedness.
- **KG structural bias:** Underrepresented edge types and highly connected pathway nodes may dominate message passing and attribution signals.

## Future Directions

- **Adaptive KG refinement:** Filter and simplify KG, optimize architecture, link to interpretability, then retrain embeddings to compare performance.
- **Temporal modeling:** account for data temporality.
- **Interpretability convergence:** assess the interpretability approaches in more depth, reapply to optimized KG, and check convergent patterns.
- **Advanced visualization:** Highlight identified biomarkers and interactive exploration of insights.