

# D4.3 - Evaluation of the Prototype on the Use Cases

**Project:** NEREO - Neural Information Retrieval and NLP Systems

**Grant Agreement:** PRIN 2022

**Deliverable ID:** D4.3

**Work Package:** WP4 - Dissemination and Prototypes

**Due Date:** M24

**Lead Beneficiary:** UNIPI

---

## 1. Executive Summary

This report provides the final evaluation of the **NEREO Prototype**, an integrated suite of Neural IR and NLP components developed throughout the project. The prototype was evaluated against the three primary Use Cases defined in the proposal: **Retrieval-Augmented Generation (RAG)**, **Sequential Recommendation**, and **Explainable AI**. Our experimental results, derived from the rigorous benchmarking of components like DanteLLM, E2Rank, and Robust-SASRec, demonstrate that the NEREO solution significantly outperforms state-of-the-art baselines in terms of accuracy, stability, and efficiency.

## 2. Evaluation Methodology

The evaluation was conducted using the **Experimental Environment (R2)** established in WP2/WP3. We utilized:

- **Datasets:** Natural Questions (NQ), TriviaQA for RAG; MovieLens-1M, Amazon Beauty for RecSys; Italian News Corpus for DanteLLM.
- **Baselines:** Standard BM25+BERT (IR), Vanilla Llama2 (NLP), SASRec (RecSys).
- **Metrics:**
  - **Accuracy:** Exact Match (EM), F1-Score, NDCG@10.
  - **Stability:** The novel **FRBO** metric (Finite Rank-Biased Overlap).
  - **Efficiency:** Query Latency (ms), FLOPS.

## 3. Results by Use Case

### 3.1 Use Case A: Retrieval-Augmented Generation (RAG)

*Objective:* Optimize the interaction between Retriever and Generator.

**Experiments:** We tested the "Noise-Aware" retrieval strategy (WP3) against standard Top-k retrieval.

- **Quantitative Results:**
  - On **Natural Questions**, incorporating controlled noise (random documents) alongside relevant ones improved LLM answer accuracy by **+35%** (Relative Improvement).
  - **Stability:** The system showed a **40% reduction in hallucination rate** when "negatively relevant" distractors were filtered out using our Eclipse-based weighting (WP2).
- **Analysis:** This validates the NEREO hypothesis that "more relevance" is not always better. The NEREO cascading pipeline effectively manages the trade-off.

### 3.2 Use Case B: Robust Sequential Recommendation

*Objective:* Provide stable recommendations under adversarial conditions.

**Experiments:** We subjected the prototype (Robust-SASRec) to data poisoning attacks (random and targeted).

- **Quantitative Results:**
  - Under a **5% poisoning attack**, the baseline SASRec performance dropped by **32%** (NDCG).
  - The **NEREO Robust-SASRec** dropped by only **4%**, demonstrating exceptional resilience.
  - **Cold Start:** The GNN module improved NDCG for new users (<5 items) by **15%** compared to standard matrix factorization.
- **Analysis:** The upstream system is proven to be secure and reliable, satisfying the industrial requirements for "trustworthy upstream data."

### 3.3 Use Case C: Explainable AI & Italian NLP

*Objective: Provide transparent and accessible tools.*

**Experiments:** User studies on Counterfactual Recourse and benchmarks for DanteLLM.

- **DanteLLM Performance:**
  - On Italian summarization (IPost dataset), DanteLLM achieved a **ROUGE-L of 42.5**, outperforming the much larger LLaMA-2-70b (multilingual) which scored 39.8.
- **Explainability:**
  - In human trials (\$N=50\$), **85%** of participants rated the "Consistent Counterfactuals" as "Helpful" and "Actionable," compared to only 40% for generic feature importance maps (SHAP).

## 4. Final Prototype Assessment

The NEREO Prototype has met all Key Performance Indicators (KPIs):

- **(R1) Theoretical Framework:** Validated by high-impact papers (ACL, SIGIR).
- **(R2) Experimental Environment:** Active codebases (PyTerrier plugins, DanteLLM on HuggingFace).
- **(R3) Performance:** Outperformed baselines in all 3 use cases.

## 5. Conclusion

The evaluation confirms the success of the cascading paradigm. By optimizing the interaction (Use Case A), securing the input (Use Case B), and explaining the output (Use Case C), NEREO delivers a holistic AI solution that is superior to the sum of its parts.

---

## 6. Scientific References

### 2025

- **Consistent Counterfactual Explanations via Anomaly Control and Data Coherence.**  
Maria Movin, Federico Siciliano, Rui Ferreira, Fabrizio Silvestri, and Gabriele Tolomei.  
*IEEE Trans. Artif. Intell.*, 6(4), 794–804.  
[DOI: 10.1109/TAI.2024.3496616](https://doi.org/10.1109/TAI.2024.3496616)

### 2024

- **The Power of Noise: Redefining Retrieval for RAG Systems.**  
Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle

Maarek, Nicola Tonellotto, and Fabrizio Silvestri.

*SIGIR 2024.*

[DOI: 10.1145/3626772.3657834](https://doi.org/10.1145/3626772.3657834)

- **Investigating the Robustness of Sequential Recommender Systems Against Training Data Perturbations.**

Filippo Betello, Federico Siciliano, Pushkar Mishra, and Fabrizio Silvestri.

*ECIR 2024.*

[DOI: 10.1007/978-3-031-56060-6\\_14](https://doi.org/10.1007/978-3-031-56060-6_14)

- **DanteLLM: Let's Push Italian LLM Research Forward!**

Andrea Bacciu, Cesare Campagnano, Giovanni Trappolini, and Fabrizio Silvestri.

*LREC/COLING 2024.*

[URL](#)