



**H.F.R.I.**  
Hellenic Foundation for  
Research & Innovation

**Description of the funded research project**  
**2nd Call for H.F.R.I. Research Projects**  
**to Support Post-Doctoral Researchers**

**Title of the research project:** Responsible by Design Entity Resolution

**Principal Investigator:** Vasilis Efthymiou

**Reader-friendly title:** ResponsibleER

**Scientific Area:** Mathematics and Information Sciences

**Institution and Country:** IBM Research, USA

**Host Institution:** Foundation for Research and Technology, Hellas (FORTH)

**Collaborating Institution(s):** Tampere University, Finland

**Project webpage**  
**(if applicable):** <https://isl.ics.forth.gr/ResponsibleER/>

**Budget:** 130,000 €

**Duration:** 36 months



## Research Project Synopsis

Big data management promises to bring a significant improvement in people's lives, accelerating knowledge discovery, research and innovation. However, in the last few years, there is an increasing concern regarding the lack of *fairness*, *diversity*, and *transparency* of data-driven algorithms supporting decision-making, raising a call for responsible by design automated decision-making systems. So far, efforts for responsible decision making have mostly focused on Machine Learning algorithms, assuming that they have been trained on high-quality data, ignoring the underlying complex pipelines that may have produced such data. A core data pipeline for producing such data is entity resolution (ER), which discovers and unifies descriptions that correspond to the same real-world entities.

In this project, we target ER systems that are responsible by design, in particular when decisions about which entity descriptions should be resolved first need to be made with respect to a given budget. The objectives of ResponsibleER are: (a) to enrich the diversity of resolved entities, (b) to ensure fairness of resolved entities, and (c) to enhance the transparency of ER systems. For (a), we are interested in formalizing progressive ER as an optimization problem, with the objective of maximizing the diversity. For (b), we are interested in measures of centrality for matching candidates in an entity graph processed by a progressive ER algorithm, then ensuring that all groups are fairly represented in the results. For (c), we need to extend the indices used for (a) and (b) and provide meaningful explanations regarding the intermediate decisions taken throughout an ER process (e.g., indexing, matching). While all three problems have been defined as major challenges recently by EU and US regulators, to the best of our knowledge, there is no other work in ER that has ever studied any of those objectives.

## Project originality

Existing ER systems aim to maximize the number of true matches by favoring the resolution of highly similar, over nearly similar descriptions, ignoring the quality of analytics conducted on the resulting datasets. Those issues become prevalent in progressive ER, where only a subset of the input entities, those expected to maximize a target benefit (e.g., number of true matches), can be selected to be resolved in order to respect a given budget constraint. We argue that a responsible data analysis requires to ensure diversity and fairness when resolving an entity graph, as well as transparency of core ER processing modules, such as entity indexing and matching. The main innovation of ResponsibleER is the matching of nearly similar descriptions, allowing to better control the depth (i.e., entity paths) and the breadth (i.e., entity attributes) of the knowledge encoded in the resolved entity graphs.

In this work, we are interested in eliminating bias regarding matching entities in conjunction with a progressive ER algorithm. To the best of our knowledge, there is no work in ER targeting fairness by design. We also attempt to bridge the gap between the optimization problem of finding the matching entities and the qualitative characteristics of the entity graph resulting after merging them. Finally, we focus on providing explicit causal explanations of the ER decisions with high conciseness. Our goal is to incorporate explainable planning for the scheduling phase of progressive ER, in which the choice of the comparison schedule (i.e., which comparison to perform next) should be properly justified over alternative schedules.

## Expected results & Research Project Impact

Despite the societal benefits of today's data-driven world, EU and US regulators, as well as policy and legal scholars argue that companies that collect and use our data have a responsibility to ensure equitable user treatment. Clearly, unfair or discriminatory effects in the outcome of analytical tasks can offend or even harm users, and cause mistrust, and potentially loss of revenue. Recent findings that some AI systems may reinforce discrimination have boosted regulations like the EU's "Right to Explanation".

ResponsibleER aims to interpret the inner representations and decisions of black-box models related to data matching tasks and lay the foundations for ER systems that can be audited and "tested for fairness" by design. The expected impact will affect the scientific community in the short and long term, but will also open opportunities for societal and economic impact in the long term. Our vision is to draw the public's and more importantly the researchers' attention to the potential societal risks that ER may bring and set the ground for a new generation of responsible by design ER systems that anyone can trust.

Regarding the scientific impact, ResponsibleER is an ambitious project in many aspects, since it will:

- introduce diversity-driven progressive ER, a problem that has never been studied before.
- be the first work in ER targeting group fairness by design.
- be the first work in ER that explicitly provides causal explanations by design.
- set the ground for responsible by design ER and constitute the first baseline for future systems.

As the outcomes of our project turn into mature software tools, they will gradually offer potential for commercial exploitation. Our prototype implementation will be available as open source software, licensed for commercial use so that enterprises can use them to design and develop responsible ER tools and applications.

## The importance of this funding

As the PI of this project, the funding of ResponsibleER has a profound impact on my career for two main reasons.

First, it sets the ground for research that is targeting to change the current focus of ER and embrace a responsible aspect that was missing from the field, leading to serious social repercussions. As researchers, it is inspiring to work for projects that we believe will make a positive change in the (digital) world around us.

Second, the funding of this project, gave me the opportunity to re-establish in Greece at a turning point in my career, after a postdoc at IBM Research - Almaden, USA, at which point the alternative was to continue my career abroad (perhaps permanently). The opportunity given by HFRI to continue my research in Greece as the PI of this research project and acquire new knowledge, experience and skills was of utmost importance.



**H.F.R.I.**  
Hellenic Foundation for  
Research & Innovation

## COMMUNICATION

185 Syggrou Ave. & 2 Sardeon St. 2  
171 21, N. Smyrni, Greece  
+30 210 64 12 410, 420  
communication@elidek.gr  
[www.elidek.gr](http://www.elidek.gr)