

Identifying Biologically Compatible Experiments via Semantic Similarity

Núria Fàbrega¹, Ian Simpson¹, Kenneth Baillie², Mirella Lapata¹

1. School of Informatics, The University of Edinburgh, Edinburgh, United Kingdom,
2. Baillie Gifford Pandemic Science Hub, University of Edinburgh, Edinburgh, UK

There are **hundreds of thousands of transcriptomics experiments publicly available** in repositories such as GEO, with substantial potential for the **validation** of results, **meta-analysis**, and **discovery** of new biological insights. Yet these datasets are **rarely reused** after the original publication, as the **metadata** is often **incomplete and unstandardised**, making it **very time-consuming** for researchers to determine which studies are biologically compatible, leaving this **golden source of information mostly untapped**.

Can we identify biologically compatible experiments by comparing their metadata?

We build a pipeline that enriches experiment metadata and measures semantic similarity to retrieve compatible studies

