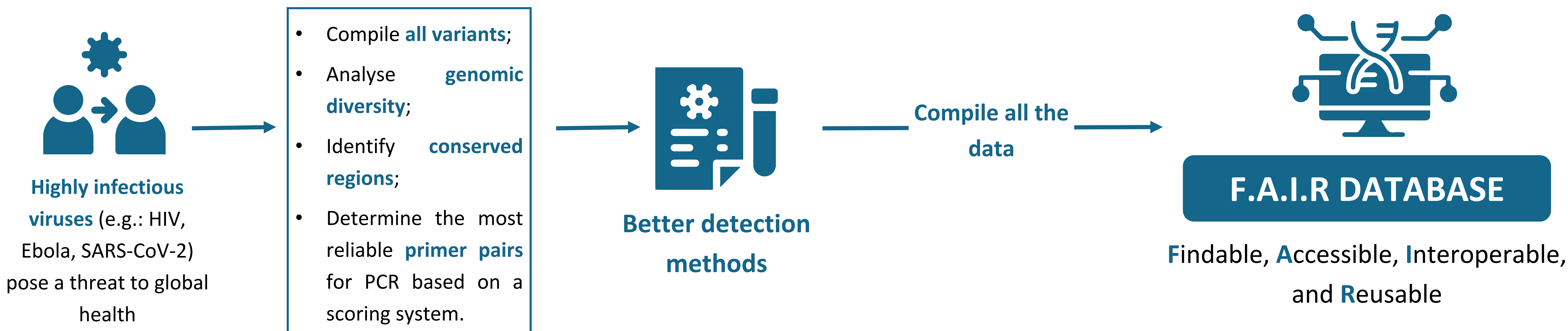


INTRODUCTION

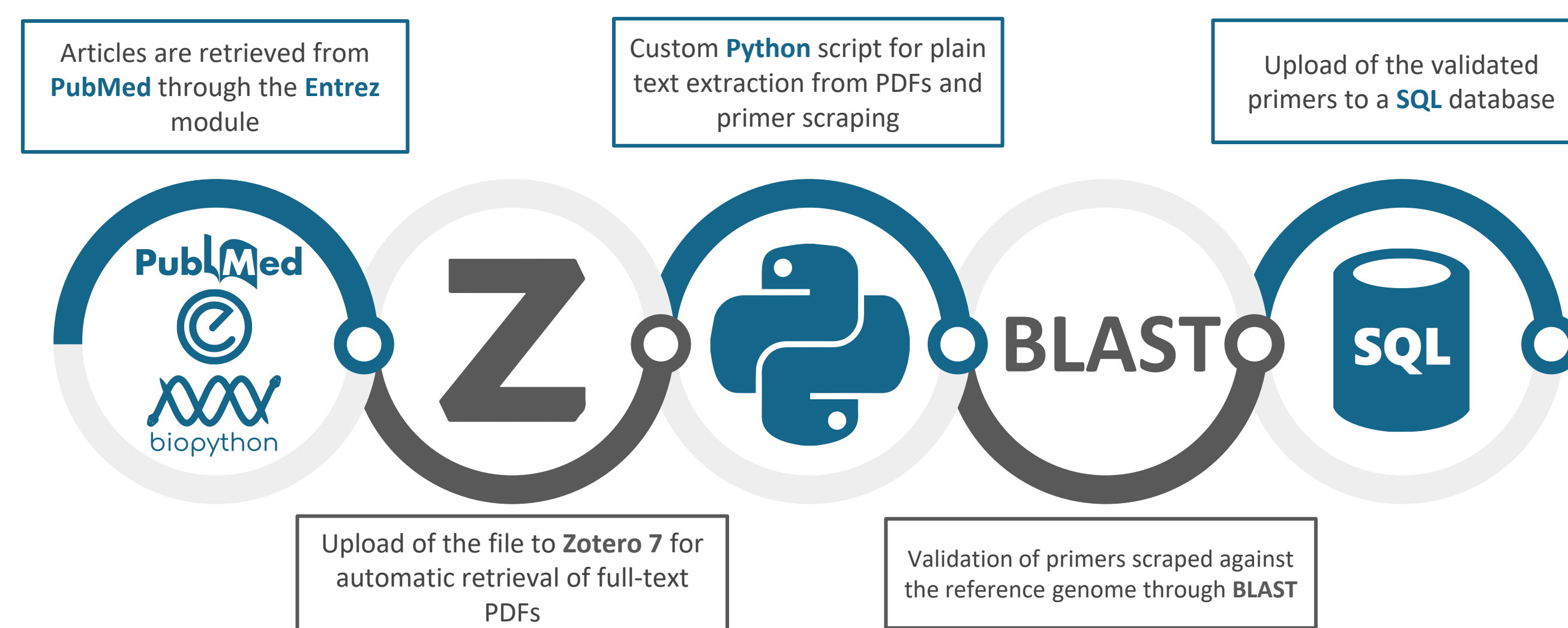
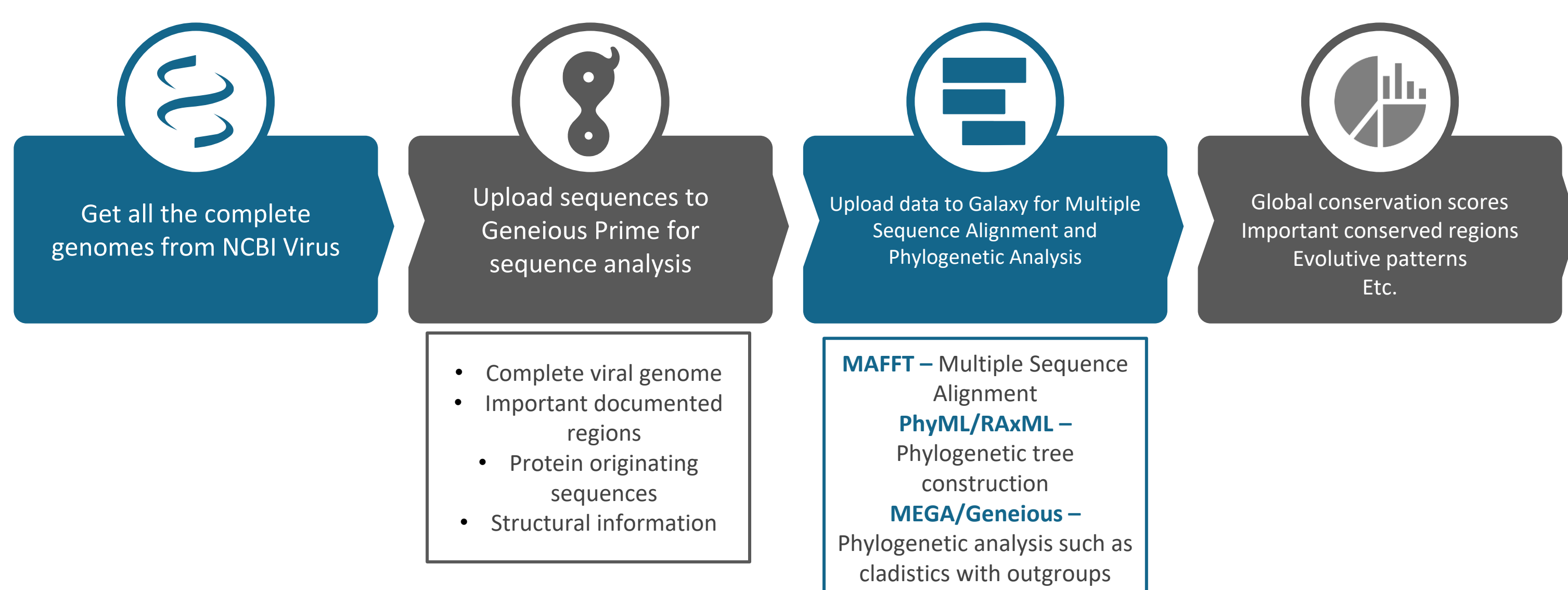
AIM



METHODOLOGY

VIRAL MULTI-OMICS DATA MINING

PRIMER RETRIEVAL



RESULTS

Retrieved:

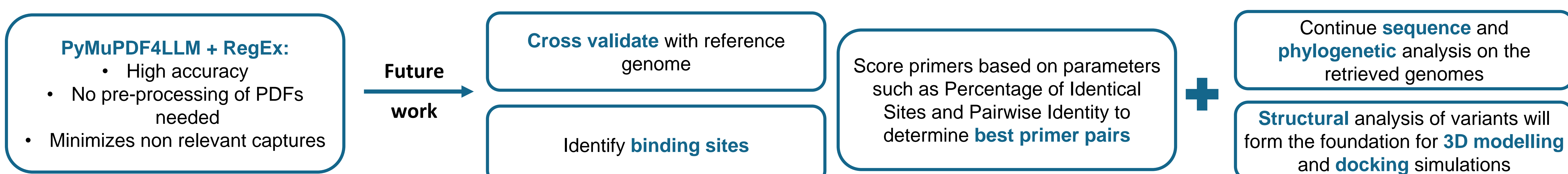
- 7304 sequences for HIV
- 658 sequences for Ebola
- 321916 sequences for SARS-CoV-2

- NotebookLM** - Large Language Model-based approach. Extraction is done using **user prompts** after uploading a batch of PDFs.
- DonutAI/OpenLaMa-7b** - Large Language Model-based approach. **DonutAI** handles **PDF to plain text** conversion and **OpenLaMa-7b** handles **primer extraction** from text.
- RegEx + PyMuPDF4LLM** - Classic approach. **PyMuPDF4LLM** handles **PDF to plain text** conversion and **primers** are captured using **regular expressions**.

Table 1. Results of the primer scraping process using three different models.

Model	Accuracy (%)	Direct export to SQL?	Execution time (s)
RegEx + PyMuPDF4LLM	71,43	Yes	41.8
DonutAI/OpenLaMa-7b	0	No	300+
NotebookLM	38,84	No	67.2

FINAL CONSIDERATIONS AND FUTURE WORK



Acknowledgments:

This research was supported by Portuguese national funds through the Foundation for Science and Technology (FCT) within the scope of UIDB/04423/2020 (CIIMAR) and UIDP/04423/2020 (CIIMAR), 10.54499/LA/P/0008/2020,10.54499/UIDP/50006/2020, UIDB/04050/2020 (UMinho CBMA - <https://doi.org/10.54499/UIDB/04050/2020>), 10.54499/UIDB/50006/2020 (LAVQ), UIDB/00319/2020 (ALGORITMI/LASI) and UIDB/00127/2020 (IEETA/LASI - doi.org/10.54499/UIDB/00127/2020). JC acknowledges the FCT funding for his research contract established under the transitional rule of Decree Law 57/2016, amended by Law 57/2017. D.P. is funded by national funds through FCT-Fundação para a Ciência e a Tecnologia, I.P., under the Scientific Employment Stimulus—Institutional Call—reference CEECINST/00026/2018.