## Unidentified dataset dependencies may inflate machine learning performance metrics in geochemistry studies

TARYN SCHARF<sup>1</sup>, MATTHEW DAGGITT<sup>2</sup>, LUC DOUCET<sup>3</sup> AND CHRISTOPHER L. KIRKLAND<sup>1</sup>

<sup>1</sup>Timescales of Mineral Systems Group, Curtin Frontier Institute for Geoscience Solutions, School of Earth and Planetary Sciences, Curtin University, Perth, WA 6103, Australia <sup>2</sup>University of Western Australia <sup>3</sup>Curtin University

Machine learning algorithms have been widely used to develop predictive models on *big data* in the geosciences. The reliability of reported model performance metrics is important if models are to inform decision making. However, model performance risks being overstated when unidentified data dependencies are present in geological datasets. Data dependencies can result in information leakage between training and testing subsets, causing a model to have artificially higher performance on the testing subset and reduced generalisability. Overstated model performance may contribute to the scientific 'reproducibility crisis' and risks decreasing confidence in machine learning applications.

We review recent peer-reviewed geochemical publications, assessing what proportion of the described methodologies explicitly consider data-leakage and what proportion may be vulnerable to overly optimistic results. Geochemical datasets often comprise multiple analyses drawn from a single sample, which imposes a hierarchical structure on the dataset. We demonstrate the effects of hierarchical data structures on algorithm performance by simulating a dataset of samples comprising multiple analyses that are related through latent variables (e.g., individual analyses linked to a common compositional sample). Simulation studies provide insights into the conditions that exacerbate the impact of data leakage in hierarchical datasets. Finally, we move beyond simulation studies by accessing the magnitude of the effect of information leakage within a subset of published datasets.

This work highlights that dependencies are often present in geological datasets. If these dependencies are unidentified and untreated, they may result in over-optimistic estimates of model performance.