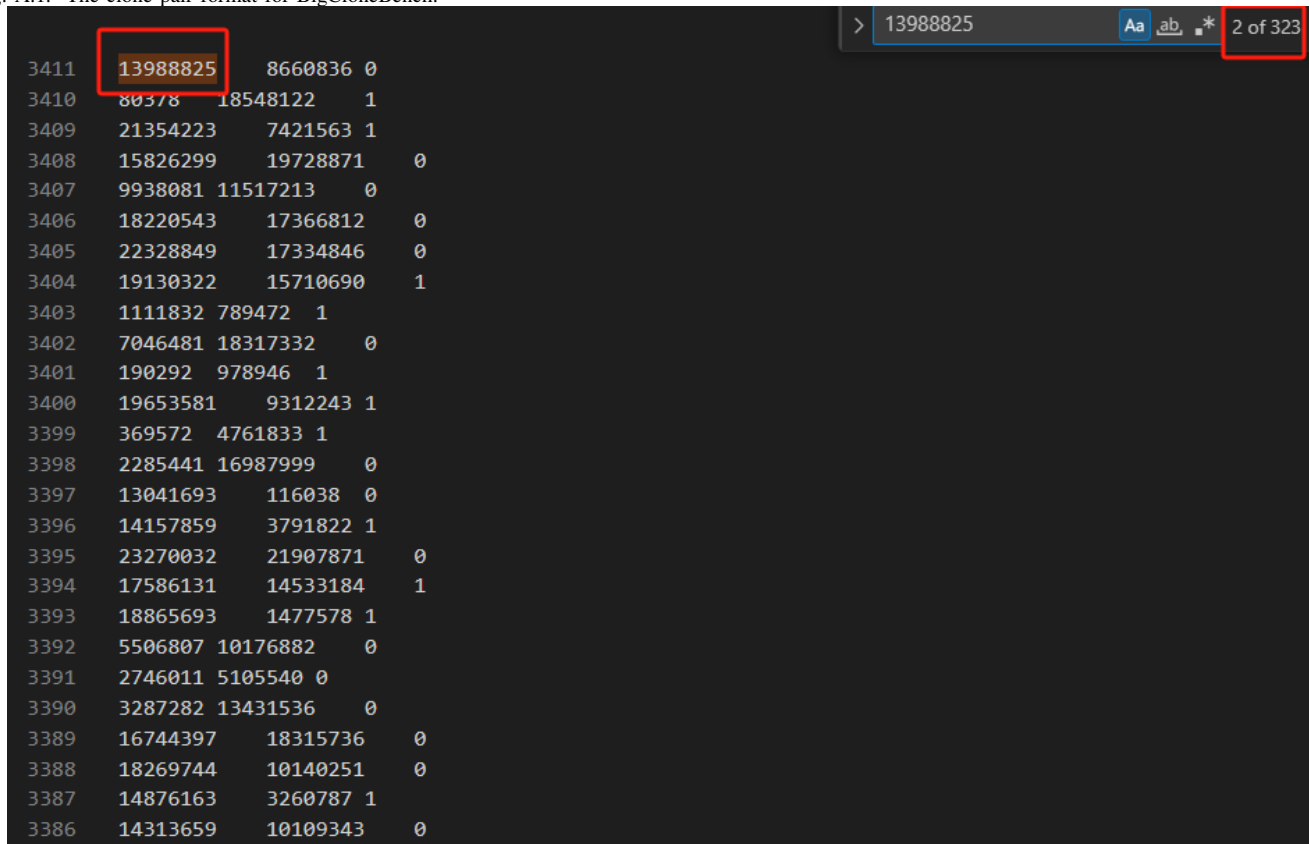# Appendix: TransformCode: A Contrastive Learning Framework for Code Embedding via Subtree Transformation

Zixiang Xian, Rubing Huang, Dave Towey, Chunrong Fang, Zhenyu Chen

Figure A.1 shows how the pairs were formatted in the BigCloneBench dataset for unsupervised code-clone detection. We used 7302 unique training samples from the BigCloneBench dataset, which were formatted in pairs to include both clones and non-clones. Our TransformCode framework was trained on the unique samples, and the same dataset settings were applied to other unsupervised tasks, maintaining consistency with the BigCloneBench dataset.

Fig. A.1. The clone pair format for BigCloneBench.