

Figure 3: The average Micro F1-score of AL models with acquisition size 100 on BERT, which were run with 5 different random seeds on various datasets.

across all datasets except for Yahoo (health). In addition, MMC and SHRL showed similar performance, as both methods focused on measuring the expected loss of the model. It is worth highlighting that SHRL showcased better results than MMC specifically in the RCV1 and TMC2007 datasets, this could be due to SHRL being the improvement of MMC by introducing the soft Hamming loss (?) as an alternative to the mean of the expected loss used in MMC. GPB2M, the most recent work in ML AL, demonstrated superior performance in the Yahoo (health) and Delicious datasets, highlighting its efficacy in those specific contexts. However, the robustness of this AL strategy raises a significant concern, as its performance exhibits variability not only across diverse datasets but also among various model architectures. For a more comprehensive analysis, including additional metrics such as macro-F1, precision and recall, please refer to the Appendix.

### Ablation Study

To assess the effectiveness and generalizability of BESRA, we conducted a comprehensive evaluation of three well-established neural network architectures commonly used in MLTC tasks: TextCNN, TextRNN, and BERT. These architectures have played a significant role in multilabel learn-

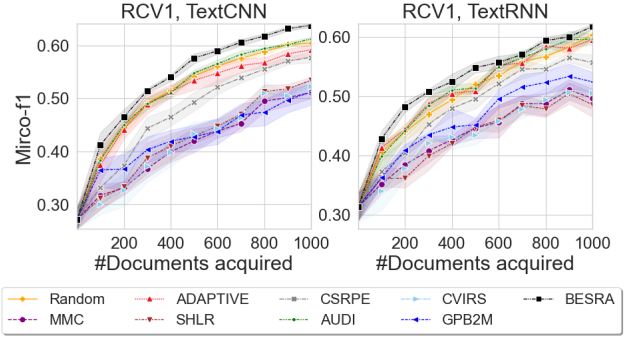


Figure 4: The average Micro F1-score of AL models with acquisition size 100 on TextCNN and TextRNN, which were run with 5 different random seeds on RCV1.

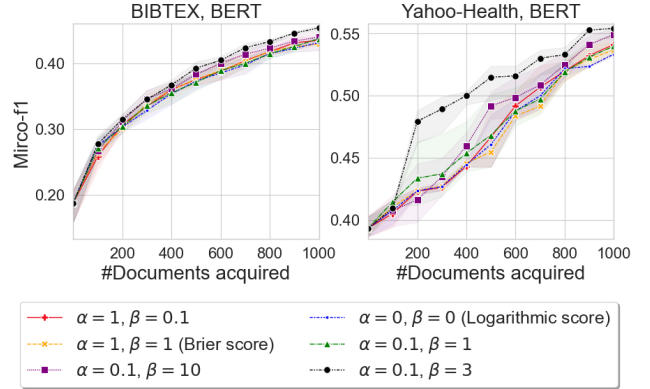


Figure 5: The average Micro F1-score of AL models with acquisition size 100 on BERT, which were run with 5 different random seeds on for Bibtex and Yahoo (health).

ing. Additionally, we investigated the impact of different Alpha and Beta values on the model’s performance. Our study aimed to answer two key research questions: (i) Does the performance advantage offered by BESRA generalise across diverse models and architectures? (ii) How do the Alpha and Beta values influence the scoring mechanism and subsequently affect the model’s performance across real-world datasets with the varied imbalance level? By addressing these questions, we aimed to gain insights of BESRA and its potential application in various MLTC tasks.

### Model Generalizability

The evaluation of BESRA is being conducted in two other models to test its generalizability. As depicted in Figure 4, BESRA consistently outperforms other baselines across TextCNN and TextRNN tested on the RCV1 dataset. Additionally, BESRA demonstrates exceptional performance compared to these models on five additional datasets included in the Appendix. These results validate that BESRA demonstrates effectiveness across a diverse range of pre-trained language models, highlighting its applicability irrespective of the specific model architecture applied.

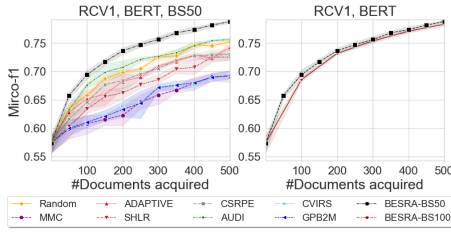


Figure 6: Left: Learning curves of ALs with batch size 50 on RCV1. Right: Learning curves for BESRA with batch sizes  $B \in \{50, 100\}$  on RCV1. All results were run with 5 different random seeds

## Batch Size

Figure 6 indicate that BESRA consistently outperforms other ALs when using a batch size of 50. Additionally, during the early stages of acquisition, BESRA performs more effectively with smaller batch sizes (50) than with larger batch sizes (100), aligning with the results reported in (??). We hypothesise that at the early acquisition stages, due to the lack of knowledge, the model trained with a limited number of samples has high uncertainty and lacks calibration, thus acquiring a large batch of samples can result in noise. Additional experiments conducted on other MLTDs can be found in the Appendix.

## Beta Parameters

In this subsection, we delve into the implications of varying Beta scores on the performance of AL, with a particular emphasis on how these scores influence the penalization behaviour of the active learner. We evaluated a range of  $\alpha$  and  $\beta$  values, associated with distinct Beta scores as shown in Eq (5). Our evaluations spanned three representative datasets, ranging from relatively balanced (i.e., BIBTEX) to highly imbalanced (i.e., Yahoo). To gain a comprehensive understanding of how different  $\alpha$  and  $\beta$  values impact the outcomes, we considered several scoring methods including the Brier score and Logarithmic score (which provide equal penalization), alongside four distinct scenarios of Beta score. These scenarios include (1) mild penalization on False Positives (FPs) where  $\alpha = 1$  and  $\beta = 0.1$ , (2) light penalization on False Negatives (FNs) with  $\alpha = 0.1$  and  $\beta = 1$ , (3) moderate penalization on FNs defined by  $\alpha = 0.1$  and  $\beta = 3$ , and finally (4) stringent penalization on FNs denoted by  $\alpha = 0.1$  and  $\beta = 10$ . Insights from Figure 5 reveal that the Brier and Logarithmic scores, as well as the scenario with  $\alpha = 1$  and  $\beta = 0.1$ , consistently underperform across various imbalance levels within MLTDs. Such results are consistent with our prior expectations, considering the inherent label imbalance challenges. A significant reason for this underperformance is the equal penalization rendered to both FN and FP outcomes by the Brier and Logarithmic scores. Additionally, the specific scenario of  $\alpha = 1$  and  $\beta = 0.1$  tends to disproportionately penalize FPs, thereby dampening performance. When evaluating the effect of  $\alpha$  and  $\beta$  values focused on penalizing FNs, we note that while light penalization settings (i.e.,  $\alpha = 0.1$ ,  $\beta = 1$ ) have a negligible impact on enhancing the active

learner’s effectiveness, the more stringent configuration of  $\alpha = 0.1$  and  $\beta = 10$  offers notable improvements, especially in highly imbalanced datasets. However it does not necessarily culminate in the optimal active learner. Instead, a moderate penalization strategy with  $\alpha = 0.1$  and  $\beta = 3$  consistently stands out as the most effective across MLTDs.

## Conclusion

We have introduced BESRA, a novel acquisition strategy for MLAL. This generalizes the recently published BEMPS using the Beta family of proper scoring rules, which allow customizable asymmetric scoring rules that effectively address the challenges such as imbalanced data associated with multi-label learning. Moreover, by our methodical construction, the use of BESRA provably converges to optimal solutions. Through empirical studies conducted on synthetic and real-world datasets, we have demonstrated the effectiveness of BESRA in acquiring highly informative samples for multi-label active learning, consistently surpassing seven existing acquisition strategies. This finding highlights the crucial role of Beta Scoring Rules and their great potential for AL with tailored acquisition strategies. Future research can further explore combinations of Alpha and Beta values for specific datasets, addressing a current limitation of BESRA. Quia quas officia labore ex ad magni id voluptatem ut, quia sit odio ab at aperiam nesciunt in omnis quam ipsa porro, asperiores laudantium dignissimos beatae a commodi, ex quasi repellat modi dignissimos veritatis explicabo necessitatibus ut ipsa, quaerat quam pariatur velit voluptatum explicabo facere?Nulla minus excepturi quo provident totam consequuntur cumque nostrum facilis, laudantium sapiente praesentium architecto ipsum eius veritatis recusandae in, et nobis sed repellat eaque ducimus, quisquam quo corrupti harum officia fugit?Fuga numquam consequatur, tempore aliquid ducimus?Magni cupiditate doloremque eius adipisci aut assumenda quo corrupti modi, voluptatem commodi fugiat architecto minus voluptas, dolores iusto corrupti distinctio perspiciatis quaerat quis illo sunt vel tenetur, ad dolorem incidunt sed dolorum quos saepe cum suscipit ab quisquam ipsa. Quos dicta similique illo tempora omnis facere eveniet explicabo ipsum, similique rerum in obcaecati suscipit. Qui eius pariatur dolorem facere fugiat, corporis eos ad eum, eligendi dignissimos dolorum totam. Quisquam ipsa reiciendis nostrum modi blanditiis officia unde nobis explicabo, eius nisi ab debitis odit magni omnis tempore, aliquid aut modi iste non ducimus, aperiam sint tempora recusandae quidem dolorum distinctio minima laboriosam incidunt?Facilis sit ipsam iure fugiat perferendis cum, error illo quis est assumenda, ad quis cupiditate ut quos itaque dolorum asperiores cumque maxime incidunt. Iste beatae debitis aut earum dignissimos sed, exercitationem fugit sunt iusto iste velit sequi a voluptatum, hic sapiente repudiandae sunt quas doloremque aliquid magni corrupti cupiditate ipsa, explicabo quaerat laboriosam enim tenetur et?Quasi temporibus dolorem nobis nostrum consequatur blanditiis vero provident aperiam, minus cum officia architecto est corporis ut qui molestiae recusandae, sunt possimus illum tempora sed sequi voluptate vitae obcaecati doloremque