

**FIGURE 6** Micro average comparison**TABLE 4** Class-specific performance evaluation of different models

(a) Results of Conditional Random Fields-2 (CRF2) model.			
Class	Precision	Recall	F-measure
B-math	83.8%	65.2%	73.2%
I-math	80.4%	75.0%	77.6%
O-math	99.2%	99.8%	99.4%
Macro average	87.8%	80.0%	83.4%
Micro average	98.7%	99.1%	98.8%
(b) Results of BIMExD model.			
Class	Precision	Recall	F-measure
B-math	73.8%	67.8%	70.7%
I-math	82.5%	80.4%	81.4%
O-math	99.4%	99.6%	99.5%
Macro average	85.3%	82.6%	83.9%
Micro average	98.9%	98.8%	98.9%
(c) Results of HIMExD model.			
Class	Precision	Recall	F-measure
B-math	79.5%	72.4%	75.8%
I-math	87.9%	78.8%	83.1%
O-math	99.4%	100%	100%
Macro average	88.9%	83.6%	86.1%
Micro average	99.0%	99.1%	99.0%

Overall results and observation of hybrid model HIMExD (presented in Table 5) are similar to the BIMExD model. Possible reasoning behind this is that the model retains the benefits of the carefully selected features (as done in CRF-2) and the automated feature engineering (as done in BIMExD), and the combination leads to better performance. The feature based model CRF-2 had a poor recall for the B-Math class, which indicates that the features were not able to detect many actual B-Math instances. Since the precision of O-Math is almost 100%, some of these false negatives for B-Math are actually being detected as I-Math—which may still be Ok as the complete chain of the mathematical expressions will still be detected. Performance for O-Math class is boosted even further, and all examples for this class are successfully detected by the method, resulting in a 100% recall. CRF2 precision is higher than the precision of HIMExD for B-Math class and BIMExD recall is better than recall of HIMExD for I-Math class. However, HIMExD is better than CRF2 and BIMExD models in terms of F-measure for all the classes. Overall, the proposed neural methods (BIMExD and HIMExD) achieve significantly better performances over the baseline approaches taken from literature.