

# 000 001 002 003 004 005 BAYESIAN POST TRAINING ENHANCEMENT OF 006 REGRESSION MODELS WITH CALIBRATED RANKINGS 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026

## ABSTRACT

027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1498  
1499  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1798  
1799  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895

054 with (i) a learned temperature that calibrates the sigmoid slope of the Bradley-Terry model, and (ii)  
 055 an accuracy-aware soft gate that regulates the influence of the ranker.  
 056

057 Specifically, our contributions are:

058

- 059 1. We introduce BAYES-ECR which combines a regressor likelihood with a calibrated ranker  
 060 likelihood to enhance predictions without retraining. We prove that the state-of-the-art,  
 061 RankRefine (Wijaya et al., 2025), is a special case under the Gaussian assumption.
- 062 2. We show how an uncalibrated Bradley-Terry model can bias the estimates when ranker  
 063 likelihood dominates, and we provide a temperature calibration with accuracy-aware soft  
 064 gating mechanism to mitigate the issue.
- 065 3. Across 12 cross-domain datasets (including real-world molecular datasets), BAYES-ECR  
 066 significantly improves over existing post-training enhancement methods (Wijaya et al.  
 067 (2025); Yan et al. (2024), and 3 other baselines) and remains effective when using imperfect  
 068 LLM rankers. Specifically, BAYES-ECR achieves 19.33% median MAE reduction relying  
 069 on 30 reference samples and a ranker with 65% accuracy, which translates to a stunning  
 070 97.65% relative improvement compared to RankRefine’s 9.78% median MAE reduction.

071 Together, these results position BAYES-ECR as a practical and principled way to leverage readily-  
 072 available pairwise information to enhance scalar regressor in data-scarce domains. Source code will  
 073 be made public upon publication.

## 074 2 BACKGROUND

075

076 **Pairwise Comparison Models.** In pairwise (binary) comparisons, the probability of item  $x_i$  is  
 077 preferred to  $x_j$  is often modeled as  $P(x_i \succ x_j) = F(y_i - y_j)$ , where  $y$  are latent scores and  $F$   
 078 is a cumulative distribution function (Cattelan, 2012). Classic pairwise comparison models include  
 079 Thurstone-Mosteller with a Gaussian *link* function (Thurstone, 1927; Mosteller, 1951) and Bradley-  
 080 Terry with a logistic *link* function (Bradley & Terry, 1952). These models provide a one-dimensional  
 081 likelihood for the unknown scalar label that can be estimated using MAP or MLE estimation.

082 **LLM-as-a-judge.** General-purpose LLMs have demonstrated strong performance on relative com-  
 083 parisons across various domains (Qin et al., 2024; Wu et al., 2024; Guo et al., 2023), with evidence  
 084 of better performance in pairwise comparison compared to score prediction (Zheng et al., 2023).  
 085 With widely available web APIs (OpenAI, 2025; Anthropic, 2025; Google, 2025), collecting com-  
 086 parisons at scale is increasingly easy, making LLMs a convenient external ranker when numeric  
 087 absolute labels are scarce.

088 **Regression Refinement using Pairwise Rankings.** Two representative approaches in this area are  
 089 the Projection method (Yan et al., 2024) and RankRefine (Wijaya et al., 2025). Projection constrains  
 090 the regressor’s prediction to a feasible interval implied by non-contradictory pairwise outcomes.  
 091 Meanwhile, RankRefine fuses the regressor’s output with a rank-only estimate via inverse variance  
 092 weighting. Our proposed method introduces a different paradigm: it reframes fusion as a Bayesian  
 093 inference, where the regressor and ranker contribute likelihoods that jointly determine the posterior.

094 **Relation to Learning from Human Feedback.** Recent works on fine-tuning with human feedback  
 095 (Christiano et al., 2017; Ziegler et al., 2019; Ouyang et al., 2022) also use pairwise rankings modeled  
 096 via Bradley-Terry, but with a different goal. They (i) learn a global reward model and (ii) fine-tune a  
 097 pretrained model accordingly. In contrast, we apply a per-query, post training prediction correction  
 098 by combining the regressor’s likelihood with a calibrated ranker likelihood. This produces a one-  
 099 dimensional posterior without modifying the regressor’s parameters.

## 100 3 METHOD

101 In this section, we describe the base formulation of BAYES-ECR (Section 3.1), along with the  
 102 analyses of its behavior (Section 3.2) and proposed modifications for improving the performance  
 103 (Section 3.3). We provide the detailed proofs of the analysis in the Appendix.

108 3.1 BAYESIAN INFERENCE ON EXPERT RANKINGS TO IMPROVE REGRESSION MODELS  
109

110 Let  $f$  be a regressor trained on  $\mathbb{C} = \{(x_j, y_j)\}_{j=1}^N$ , with scalar labels  $y_j \in \mathbb{R}$  and prediction  $f(x_j) = \hat{y}_j^{\text{re}}$ . Let  $R$  be an expert pairwise ranker that returns a binary comparison:  $R(x_a, x_b) = 1$  if it predicts  $y_a > y_b$ , and  $R(x_a, x_b) = 0$  otherwise. We are given a reference set  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^k$  with known  $y_i$ , which can be, but not always is, the regressor's training set  $\mathbb{C}$ . For a query  $x_0$  with unknown  $y_0 \in \mathbb{R}$ , we collect *expert rankings* against all references:  $\mathbb{G} = \{r_i = R(x_0, x_i)\}_{i=1}^k$  where  $r_i \in \{0, 1\}$ .

116 Assuming conditional independence of regressor and ranker given the true  $y_0$ , the Bayes' rule yields  
117

$$\underbrace{p(y_0 \mid \hat{y}_0^{\text{re}}, \mathbb{G})}_{\text{posterior}} \propto \underbrace{p(\hat{y}_0^{\text{re}} \mid y_0)}_{\text{reg. likelihood}} \underbrace{p(\mathbb{G} \mid y_0)}_{\text{rank likelihood}} \underbrace{p(y_0)}_{\text{prior}} \quad (1)$$

121 with Gaussian regressor likelihood  $p(\hat{y}_0^{\text{re}} \mid y_0) = \mathcal{N}(\hat{y}_0^{\text{re}}; y_0, \sigma_{\text{re}}^2)$  and rank likelihood

$$p(r_i = 1 \mid y_0, y_i) = s(y_0 - y_i), \quad s(z) = 1/(1 + e^{-z}), \quad (2)$$

$$p(\mathbb{G} \mid y_0) = \prod_{i=1}^k s(y_0 - y_i)^{r_i} (1 - s(y_0 - y_i))^{1-r_i}, \quad (3)$$

124 which is the product of all pairwise comparisons in  $\mathbb{G}$  under the Bradley-Terry model (Bradley &  
125 Terry, 1952).

126 With the two likelihoods, the *base* formulation of BAYES-ECR enhances a regressor prediction with  
127 expert rankings using a maximum *a posteriori* (MAP) estimation (Murphy, 2022) which maximizes  
128  $y$  with the objective function  
129

$$\mathcal{L}(y) = -\frac{1}{2\sigma_{\text{re}}^2} (\hat{y}_0^{\text{re}} - y)^2 + \sum_{i=1}^k \left[ r_i \log s(y - y_i) + (1 - r_i) \log(1 - s(y - y_i)) \right] + \log p(y).$$

(4)

130 If we assume a flat (uninformative) prior, i.e.,  $p(y) = 1$ , BAYES-ECR reduces to a maximum  
131 likelihood estimation (MLE) (Murphy, 2022). The objective function of the BAYES-ECR (Equation  
132 4) is strictly log-concave so we can obtain the maximum efficiently.

133 **Lemma 3.1. (Strict log-concavity.)** If  $\sigma_{\text{re}}^2 > 0$  and  $\log p(y_0)$  is concave (including the flat prior),  
134 then  $\log p(y_0 \mid \hat{y}_0^{\text{re}}, \mathbb{G})$  is strictly concave in  $y_0$ .  $\square$

135 **Corollary 3.2. (Existence and uniqueness.)** The maximum *a posteriori* (and maximum likelihood  
136 under a flat prior) estimate exists and is unique.  $\square$

137 **Proposition 3.3. (RankRefine (Wijaya et al., 2025) is a special case of BAYES-ECR under Gaussian assumption.)** Let  $\mathcal{L}_{BT}(y)$  denote the rank-only log-likelihood and  $\hat{y}_0^{\text{ra}} = \arg \max_y \mathcal{L}_{BT}(y)$ . By assuming Gaussianity on the ranker likelihood (similar to Wijaya et al.  
138 (2025)), the second-order expansion of  $\mathcal{L}_{BT}(y)$  around  $\hat{y}_0^{\text{ra}}$  yields  
139

$$p(\mathbb{G} \mid y) \approx \mathcal{N}(y; \hat{y}_0^{\text{ra}}, \sigma_{\text{ra}}^2), \quad \sigma_{\text{ra}}^2 = \left[ \sum_{i=1}^k s(\hat{y}_0^{\text{ra}} - y_i)(1 - s(\hat{y}_0^{\text{ra}} - y_i)) \right]^{-1}. \quad (5)$$

140 Combining this with the Gaussian regressor likelihood and a flat prior gives the inverse-variance  
141 weighted (IVW) (Cochran & Carroll, 1953) estimator, i.e., RankRefine (Wijaya et al., 2025), for  
142 which the estimate  $\hat{y}_0^{\text{rr}}$  is,  
143

$$\hat{y}_0^{\text{rr}} = \frac{\hat{y}_0^{\text{re}}/\sigma_{\text{re}}^2 + \hat{y}_0^{\text{ra}}/\sigma_{\text{ra}}^2}{1/\sigma_{\text{re}}^2 + 1/\sigma_{\text{ra}}^2}. \quad \square \quad (6)$$

144 *Proof.* The gradient and curvature of  $\mathcal{L}_{BT}(y)$  are  
145

$$\mathcal{L}'_{BT}(y) = \sum_{i=1}^k (r_i - s(y - y_i)), \quad \mathcal{L}''_{BT}(y) = - \sum_{i=1}^k (s(y - y_i)(1 - s(y - y_i))). \quad (7)$$

162 The second-order Taylor expansion of  $\mathcal{L}_{BT}$  around  $y = \hat{y}_0^{\text{ra}}$  is  
 163

$$\mathcal{L}_{BT}(y) \approx \mathcal{L}_{BT}(\hat{y}_0^{\text{ra}}) + \mathcal{L}'_{BT}(\hat{y}_0^{\text{ra}}) \cdot (y - \hat{y}_0^{\text{ra}}) + \frac{1}{2} \mathcal{L}''_{BT}(\hat{y}_0^{\text{ra}}) \cdot (y - \hat{y}_0^{\text{ra}})^2. \quad (8)$$

164 The first term is constant and the second term is zero, therefore,  
 165

$$\mathcal{L}_{BT}(y) \propto -\frac{1}{2} \sum_{i=1}^k s(\hat{y}_0^{\text{ra}} - y_i)(1 - s(\hat{y}_0^{\text{ra}} - y_i))(y - \hat{y}_0^{\text{ra}})^2. \quad (9)$$

170 Similar to RankRefine, we make a strong assumption of Gaussianity for the BT likelihood,  
 171

$$\begin{aligned} p(\mathbb{G} | y_0) &\propto \exp(\mathcal{L}_{BT}(y)) = \exp\left(-\frac{1}{2} \sum_{i=1}^k s(\hat{y}_0^{\text{ra}} - y_i)(1 - s(\hat{y}_0^{\text{ra}} - y_i))(y - \hat{y}_0^{\text{ra}})^2\right) \equiv \mathcal{N}(\mu_{ra}, \sigma_{ra}^2), \\ \mu_{ra} &= \hat{y}_0^{\text{ra}}, \quad \sigma_{ra}^2 = \left[\sum_{i=1}^k s(\hat{y}_0^{\text{ra}} - y_i)(1 - s(\hat{y}_0^{\text{ra}} - y_i))\right]^{-1}. \end{aligned} \quad (10)$$

172 The joint likelihood of the conditionally independent regressor and ranker is now a product of two  
 173 Gaussians. The MLE, i.e., MAP estimate with a flat prior, then solves a weighted least-squares  
 174 problem whose maximizer is the inverse-variance weighted average.  $\square$   
 175

### 182 3.2 ANALYSIS OF POTENTIAL PERFORMANCE DEGRADATION OF BASE BAYES-ECR

184 Under the base formulation, BAYES-ECR can degrade  
 185 when the reference size  $k$  and ranker accuracy  $a$  exceed cer-  
 186 tain thresholds. The effect is dataset-dependent and arises  
 187 from a mismatch between the Bradley-Terry model and the  
 188 behavior of real-world rankers.

189 Separating the rank likelihood from the objective in Equa-  
 190 tion 4 and maximizing it yields the rank-only MLE  $\hat{y}_0^{\text{ra}}$ . If  
 191 the ranker were perfectly accurate and the likelihood per-  
 192 fectly specified, then  $\hat{y}_0^{\text{ra}} \in (y_m, y_{m+1})$ , where  $y_m < y_0 <$   
 193  $y_{m+1}$  are the closest references and  $m = \sum_i r_i$  is the num-  
 194 ber of references ranked below  $y_0$ . In practice, the sigmoid  
 195 in the Bradley-Terry model induces a *soft- vs. hard-count*  
 196 mismatch: the rank-only target solves a soft-count equation  
 197 and need not lie in  $(y_m, y_{m+1})$ .

198 **Lemma 3.4 (Rank-only MLE may target a pseudo**  
 199 **ground truth.)** Let  $u_i(y) = s(y - y_i)$  and  $m = \sum_{i=1}^k r_i$ .  
 200 The Bradley-Terry log-likelihood  $\mathcal{L}_{BT}(y)$  is strictly con-  
 201 cave and its unique maximizer  $\hat{y}_0^{\text{ra}}$  satisfies

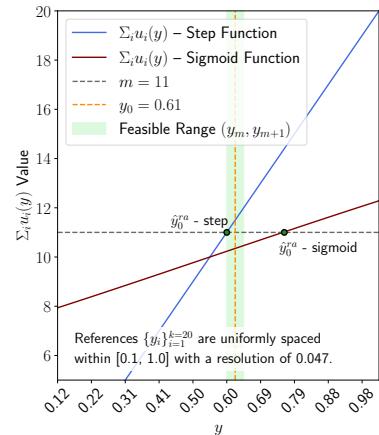
$$\mathcal{L}'_{BT}(\hat{y}_0^{\text{ra}}) = 0 \iff \sum_{i=1}^k u_i(\hat{y}_0^{\text{ra}}) = m. \quad (11)$$

202 Because  $m$  is a *hard count* while  $u_i$  are sigmoid *soft counts*,  
 203 the solution to  $\sum_{i=1}^k u_i(\hat{y}_0^{\text{ra}}) = m$  is a pseudo target  $\tilde{y}_0$ ,  
 204 which can lie outside  $(y_m, y_{m+1})$  even under a perfectly  
 205 accurate ranker, biasing  $\hat{y}_0^{\text{ra}}$ . See Figure 1 for illustration.  $\square$

206 Bias from the rank-only target can increase the overall error  
 207 if the rank likelihood dominates the regressor likelihood.

208 Rank likelihood dominance is governed by its curvature (Fisher information (Ly et al., 2017)), which  
 209 grows with  $k$ .

210 **Lemma 3.5. (Rank curvature (Fisher information) scales with  $k$ .)** The Fisher information of  
 211 the rank log-likelihood is  $I_{\text{rank}}(y) = \sum_{i=1}^k u_i(y)(1 - u_i(y))$ , so  $I_{\text{rank}}(y)$  grows linearly with  $k$ .  
 212 Moreover,  $0 \leq I_{\text{rank}}(y) \leq \frac{k}{4}$ .  $\square$



213 Figure 1: Soft- vs. hard-count mismatch from Lemma 3.4. We gen-  
 214 erate a synthetic reference set  $\{y_i\}$  and  
 215 ground truth  $y_0$ . From Equation 11,  
 216 the rank-only MLE  $\hat{y}_0^{\text{ra}}$  is the inter-  
 217 section between  $\sum u_i(y)$  and  $m$ , and  
 218  $(y_m, y_{m+1})$  forms the feasible range.  
 219 We can see that  $\sum u_i(y)$  with a step  
 220 function correctly intersects the feasible  
 221 range. However, using a sigmoid  
 222 function results in a biased intersec-  
 223 tion outside the feasible range.

216 **Lemma 3.6. (Info-weighted Newton step and dominance.)** As  $k$  grows, the optimization process  
 217 using Newton’s method (Murphy, 2022) to solve Equation 4 is dominated by the rank likelihood  
 218 term. That is, denoting  $g_{\text{tot}}(y)$  and  $I_{\text{tot}}(y)$  as the total gradient and Fisher information of both likeli-  
 219 hoods at  $y$ ,  $I_{\text{reg}}$  as the Fisher information of the regressor likelihood, and  $\tilde{y}^{\text{ra}}(y)$  as the target of the  
 220 ranker likelihood term at  $y$ , a Newton step is

$$221 \quad y \leftarrow y + \frac{g_{\text{tot}}(y)}{I_{\text{tot}}(y)} = \frac{I_{\text{reg}} \cdot \hat{y}_0^{\text{re}} + I_{\text{rank}}(y) \cdot \tilde{y}^{\text{ra}}(y)}{I_{\text{reg}} + I_{\text{rank}}(y)}, \quad (12)$$

224 which is an information-weighted average using Fisher information as the weights.  $\square$

225 Since  $I_{\text{reg}} = 1/\sigma_{\text{re}}^2$  is constant, Lemma 3.5 implies the rank likelihood term eventually dominates as  
 226  $k$  grows. Apart from Newton, for a first-order step  $y \leftarrow y + \eta g_{\text{tot}}(y)$ , the same decomposition yields  
 227  $y \leftarrow y + \eta [I_{\text{reg}}(\hat{y}_0^{\text{re}} - y) + I_{\text{rank}}(y)(\tilde{y}^{\text{ra}}(y) - y)]$ , showing similar information-weighted pull. The  
 228 dominance of the rank likelihood then follows from Lemma 3.5.

229 Dominance alone is harmless, but combined with Lemma 3.4, a biased rank likelihood target  $\tilde{y}_0$  can  
 230 steer the refinement away from  $y_0$ . Note that dominance occurs only when both  $k$  and the ranker  
 231 accuracy  $a$  are sufficiently large. Moreover, the accuracy threshold decreases with  $k$ .

232 **Lemma 3.7. (Accuracy- $k$  threshold for rank dominance.)** Define

$$234 \quad p^* = \frac{1}{k} \sum_{i=1}^k \mathbb{1}(y_0 > y_i), \quad \hat{p}(y) = \frac{1}{k} \sum_{i=1}^k u_i(y). \quad (13)$$

237 Let  $a \in [0.5, 1]$  denote the ranker accuracy. Then, at any iterate  $y$ , the rank likelihood dominates the  
 238 regressor likelihood whenever

$$239 \quad a > \frac{1}{2} + \frac{|y - \hat{y}_0^{\text{re}}|}{2k\sigma_{\text{re}}^2 |p^* - \hat{p}(y)|}. \quad \square \quad (14)$$

241 The right hand side of the inequality is the threshold accuracy  $a_{\text{thr}}(y)$ , and it decreases as  $1/k$  for a  
 242 fixed  $y$ ,  $\sigma_{\text{re}}^2$ , and  $|p^* - \hat{p}(y)| > 0$ . Therefore, the required ranker accuracy such that the rank term  
 243 dominates the Newton step shrinks as  $k$  grows.

244 **Corollary 3.8. (Rank likelihood can degrade base BAYES-ECR.)** Under the base formulation  
 245 of BAYES-ECR, the rank-only target can be biased (Lemma 3.4). This may lead to a performance  
 246 degradation when the ranker likelihood term dominates (Lemma 3.6). In general, ranker likelihood  
 247 dominance grows with the reference set size  $k$  and the ranker accuracy  $a$ . When  $k$  is large, the  
 248 required  $a$  for the ranker likelihood to dominate is lower (Lemma 3.7).

### 250 3.3 BAYES-ECR WITH CALIBRATED EXPERT RANKINGS

252 The degradation stems from the Bradley-Terry modeling mismatch (Lemma 3.4): when many pair-  
 253 wise gaps  $(y_a - y_b)$  fall on the sigmoid’s transition region rather than its saturated tails, the rank-only  
 254 target becomes biased. This is a *unit-scale* issue; datasets with small average  $|y_a - y_b|$  have higher  
 255 probability of placing many pairs on the transition slope.

256 We address this by introducing a temperature  $\tau$  to adjust the slope, i.e., replace  $u_i(y) = s(y - y_i)$   
 257 with  $v_i(y; \tau) = s((y - y_i)/\tau)$ . This (i) aligns the logistic slope with label units so that the tempered  
 258 score  $\sum_i v_i(\hat{y}_0^{\text{ra}}; \tau)$  can match the observed count  $m = \sum_i r_i$ , removing the *soft-* vs. *hard-count*  
 259 mismatch, and (ii) sets the rank curvature to

$$260 \quad I_{\text{rank}}(y; \tau) = \tau^{-2} \sum_i v_i(y; \tau) (1 - v_i(y; \tau)), \quad (15)$$

262 thereby controlling rank dominance at large  $k$ . We estimate  $\tau$  for a dataset via a one-parameter  
 263 logistic fit on the reference set, thus requiring no extra labeled data,

$$265 \quad \Pr(r = 1 | y_a, y_b; \hat{\omega}) = s(\hat{\omega}(y_a - y_b)), \quad (y_a, y_b) \in \mathbb{D}, \quad (16)$$

266 and set  $\tau = \hat{\tau}_{\text{cal}} = 1/\hat{\omega}$ . With this calibration of expert rankings, the Newton update is an  
 267 information-weighted average with temperature-controlled curvature,

$$268 \quad \boxed{269 \quad y \leftarrow \frac{I_{\text{reg}} \cdot \hat{y}_0^{\text{re}} + I_{\text{rank}}(y; \tau) \cdot \tilde{y}^{\text{ra}}(y; \tau)}{I_{\text{reg}} + I_{\text{rank}}(y; \tau)}}. \quad (17)$$

270 where  $\hat{y}^{\text{ra}}(y; \tau) = y + \frac{\sum_{i=1}^k (r_i - v_i(y; \tau))}{I_{\text{rank}}(y; \tau)}$ .  
 271

272 Setting  $\tau < 1$  increases the rank likelihood curvature (via the  $\tau^{-2}$  factor) and shifts the update  
 273 toward the rank target. When the ranker is accurate, this permits rank dominance *without* the bias  
 274 highlighted by Lemma 3.4, improving performance over the naive BAYES-ECR.

275 On the other hand, when the ranker is less accurate, overly large curvature would let noisy rank  
 276 signals dominate and harm performance. We therefore apply an *accuracy-aware soft gate*:  
 277

$$\tau(a) = 1 + (\hat{\tau}_{\text{cal}} - 1)(w(a))^\gamma, \quad \gamma \geq 1, \quad (18)$$

$$w(a) = \max(0, 2a - 1) \in [0, 1]. \quad (19)$$

281 Thus  $\tau(a) \approx 1$  when  $a \approx 0.5$ ,  $\tau(a) \rightarrow \hat{\tau}_{\text{cal}}$  as  $a \rightarrow 1$ , and intermediate  $a$  produces a smooth  
 282 interpolation. The final BAYES-ECR with calibrated expert rankings is summarized in Algorithm 1.  
 283

---

284 **Algorithm 1** BAYES-ECR Algorithm  
 285

---

286 **Inputs:**  $x_0, \hat{y}_0^{\text{re}}, \sigma_{\text{re}}^2, \mathbb{G} = \{(x_i, y_i)\}_{i=1}^k; \mathbb{D} = \{y_i\}_{i=1}^k R$ ; (optional:  $p(y)$ )

287 **Collect Comparisons:**  $r_i \leftarrow R(x_0, x_i)$

288 **Estimate Temperature:** fit  $s(\alpha(y_a - y_b))$  on labeled pairs from  $\mathbb{D}$ ; set  $\hat{\tau}_{\text{cal}} = 1/\hat{\alpha}$ . Estimate the  
 289 ranker accuracy  $a$ ; set  $\tau \leftarrow 1 + (\hat{\tau}_{\text{cal}} - 1)(\max(0, 2a - 1))^\gamma$ .

290 **MAP / MLE Optimization:** Maximize Equation 4 with  $s(\cdot)$  replaced by  $s(\cdot)/\tau$ . With Newton  
 291 steps described on Equation 17, iterate to convergence.

292 **Output:**  $\hat{y}_0^{\text{rr}}$  and approximate uncertainty  $\sigma_{\text{post}}^2 \approx (I_{\text{reg}} + I_{\text{rank}}(\hat{y}_0^{\text{rr}}; \tau))^{-1}$ .  
 293

---

294 **4 EXPERIMENTS**  
 295

296 To demonstrate the benefits of BAYES-ECR in data-scarce domains, we evaluate on 9 molecular  
 297 property prediction datasets from the TDC ADMET regression task (Huang et al., 2021): Caco-2  
 298 (Wang et al., 2016), Clearance Microsome and Clearance Hepatocyte (Di et al., 2012), log Half-Life  
 299 (Obach et al., 2008), FreeSolv (Mobley & Guthrie, 2014), Lipophilicity (Wu et al., 2018), PPBR,  
 300 Solubility (Sorkun et al., 2019), and VDss (Lombardo & Jing, 2016). For each experiment, we sample  
 301  $N = 100$  molecules from the original training split and  $L = 100$  molecules from the original  
 302 test split. We repeat this train/test resampling for 10 random seeds, similar to a Monte Carlo cross-  
 303 validation. The reference set  $\mathbb{D}$  is sampled from the training set with  $k \in \{3, 10, 20, 30, 50, 100\}$ .  
 304 Two types of base regressors are used: random forest (RF) (Ho, 1995) and multilayer perceptron  
 305 (MLP) (Rumelhart et al., 1986), both trained on the  $N$  training labels. We also use three tabular re-  
 306 gression datasets: crop-yield prediction from sensor data (Soundankar, 2025), student-performance  
 307 prediction (Cortez, 2014), international-education cost estimation (Shamim, 2025).

308 We evaluate four BAYES-ECR variants: MAP with a Gaussian prior, MLE, MLE with tempera-  
 309 ture scaling (MLE-Temp), and MLE with temperature scaling and soft gating (MLE-GatedTemp).  
 310 The full BAYES-ECR, i.e., MLE-GatedTemp, is the default variant we use in the experiments. We  
 311 compare our method to two post-training regression enhancement baselines: (i) Projection-based  
 312 approach (Yan et al., 2024) and (ii) RankRefine (Wijaya et al., 2025). Additionally, we adapt  
 313 two pairwise ranking models to create rank-only regression enhancement baselines: (i) Bradley-  
 314 Terry (rank-only MLE using the Bradley-Terry model), (ii) Thurstone (rank-only MLE using the  
 315 Thurstone-Mosteller model (Thurstone, 1927; Mosteller, 1951)). Following Wijaya et al. (2025),  
 316 we report  $\beta \equiv \text{MAE}_{\text{post}}/\text{MAE}_{\text{base}}$ , which is the ratio between the post-enhancement Mean Absolute  
 317 Error (MAE) and the MAE of the regressor-only predictions (lower is better).

318 We use two types of rankers: (i) *oracle* ranker to study the effect of ranker accuracy, and (ii) LLM  
 319 ranker to demonstrate real-world use cases. For the oracle ranker, we set  $r_i = \mathbb{1}(y_0, y_i)$  and then  
 320 flip the outcomes with probability  $1 - a$  to simulate a ranker with accuracy  $a \in [0.5, 1]$ . For  
 321 LLM rankers, we prompt publicly-available models to compare a list of pairs of molecular text  
 322 representations (SMILES (Weininger, 1988)) and measure the accuracy  $a$  on a small validation set.

323 Additional discussions and experiments on LLM usage, temperature calibration, multi-target regres-  
 324 sion, and reference set size are available in Appendix.

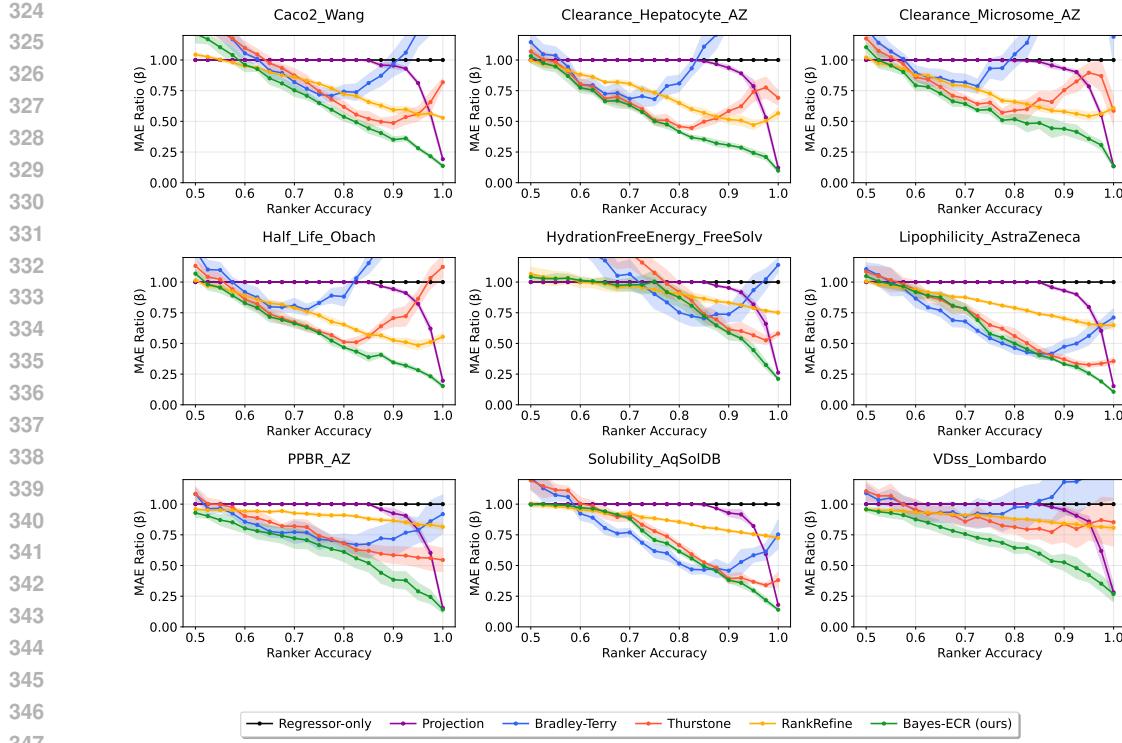


Figure 2: MAE ratio  $\beta$  (lower is better) as a function of oracle ranker accuracy  $a$  on nine TDC ADMET datasets ( $k = 30$ ). We report the mean and standard deviation across 10 random splits. BAYES-ECR (MLE-GatedTemp variant) outperforms projection and RankRefine across ranker accuracies, with especially strong gains at moderate, real-world accuracy. Results for other  $k$  values with similar trends are available in Appendix.

#### 4.1 MAIN RESULTS WITH ORACLE RANKERS

Across all nine ADMET datasets at  $k = 30$ , BAYES-ECR (MLE-GatedTemp variant) consistently outperforms previous state-of-the-art, RankRefine, and 3 other baselines (Figure 2). Gains over recent custom-built enhancement methods, RankRefine and Projection, are largest in the mid-accuracy regime (65% - 90%). At lower ranker accuracies (<65%), BAYES-ECR matches or exceeds RankRefine and clearly outperforms Projection. At high ranker accuracies, RankRefine performance plateaus, while BAYES-ECR continues to improve, converging to Projection’s performance at perfect (yet unrealistic) ranker accuracy.

#### 4.2 ABLATION: PRIOR, TEMPERATURE, AND ACCURACY-AWARE SOFT GATING

Figure 3 shows that MAP/MLE variants of BAYES-ECR can degrade as  $k$  grows, confirming that rank-dominance in the Newton step and the soft-hard count mismatch on the rank term can bias the enhancement (Corollary 3.8). The rank-only Bradley-Terry and Thurstone models in Figure 2 also exhibit similar degradation, validating that the source of degradation is the ranker likelihood. Interestingly, MAP degrades less than MLE, indicating that prior can act as a regularizer. As proposed in Lemma 3.7, the ranker accuracy threshold at which degradation appears shifts lower as  $k$  increases.

Our full method (MLE-GatedTemp variant) solves the performance degradation issue. Temperature scaling aligns the sigmoid slope to the label scale and controls the rank curvature, mitigating degradation at high ranker accuracy, but can worsen performance at low ranker accuracy (Figure 3, MLE-Temp). Adding accuracy-aware soft-gating removes this trade-off (Figure 3, MLE-GatedTemp), enabling full BAYES-ECR to deliver strong performance across the complete ranker accuracy range. Note that Figure 3 also confirms that our enhancement method works with a reference set as small as  $k = 3$ , which is highly practical in the real-world setting.

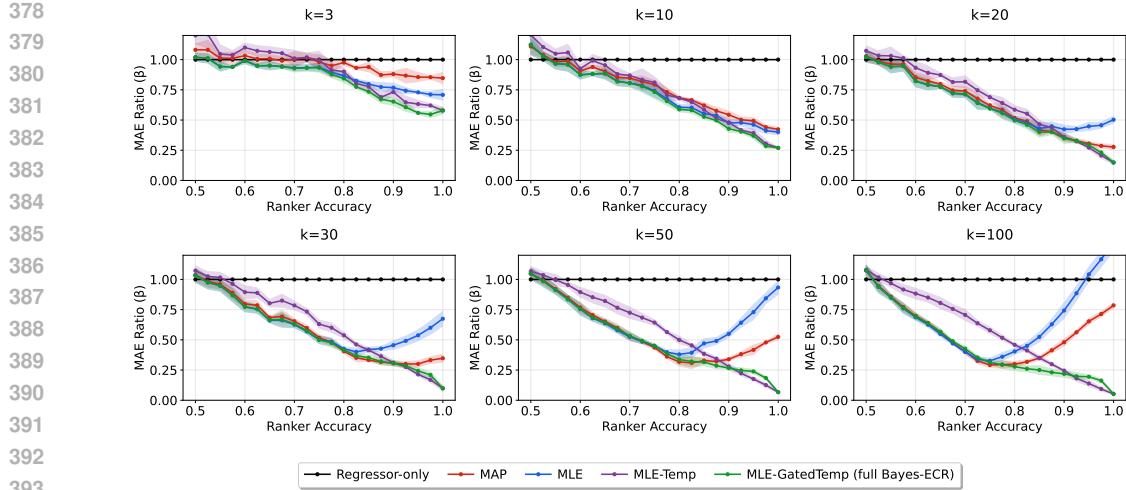


Figure 3: Effects of reference set size  $k$  and temperature scaling / soft gating on the Clearance Hepatocyte dataset. The shown methods are the four variants of BAYES-ECR, with *ours* being MLE-GatedTemp. We show the MAE ratio  $\beta$  (lower is better) as a function of oracle ranker accuracy  $a$  for  $k \in \{3, 10, 20, 30, 50, 100\}$ . Without temperature scaling, MAP and MLE variants degrade at larger  $k$  due to rank curvature dominance. Temperature (MLE-Temp) fixes the scale mismatch, and soft-gating (MLE-GatedTemp) further improves robustness when ranker accuracy is low.

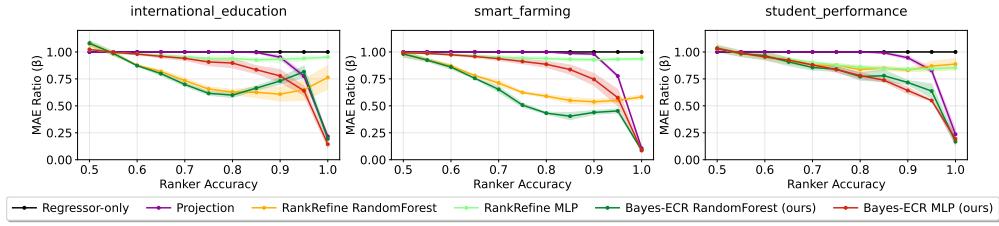


Figure 4: MAE ratio  $\beta$  on three non-molecular tabular datasets with Random Forest (RF) and Multi-layer Perceptron (MLP) as base regressor models. We show that the refinements are generalizable to other domains and regressor models. Baseline (non-refined) MAEs are: **International Education Cost**: 0.0926 (MLP), 0.0609 (RF); **Smart Farming Yield**: 0.1448 (MLP), 0.1361 (RF); **Student Exam Performance**: 0.0727 (MLP), 0.0819 (RF).

### 4.3 CROSS-DOMAIN GENERALITY AND REGRESSOR DIVERSITY

To test cross-domain and cross-model generalization, we add three non-molecular tabular datasets and run both Random Forest (RF) and Multilayer Perceptron (MLP) regression models. Figure 4 shows that BAYES-ECR yields  $\beta < 1$  across the ranker accuracy range, improving both RF and MLP models, and proving its generalization capability to other domains and regression models.

### 4.4 LLM AS IMPERFECT YET PRACTICAL RANKER

We replace the oracle ranker with off-the-shelf LLMs to obtain noisy but scalable expert rankings. With ChatGPT5 Thinking and Claude Sonnet 4 as rankers at  $N = 50$  and  $k = 20$ , BAYES-ECR generally achieves the best  $\beta$  across five ADMET datasets. Detailed results are shown in Table 1.

### 4.5 RUNTIME ANALYSIS ON VARYING REFERENCE SET SIZE

Excluding the ranking predictions, BAYES-ECR runs in less than 1 ms on an Intel i7-13700 CPU when the reference set size  $k \leq 1000$ , covering realistic scenarios in data-scarce domains. The average running time to predict expert rankings using ChatGPT5-Thinking for 1000 queries and 50

	Dataset Name	Half Life	FreeSolv	PPBR	Solubility	VDss
432	ChatGPT5 PRA (%)	62.20 $\pm$ 1.89	62.27 $\pm$ 2.00	63.29 $\pm$ 2.65	62.87 $\pm$ 3.27	65.95 $\pm$ 1.96
433	BAYES-ECR ( $\beta$ )	<b>0.952 <math>\pm</math> 0.057</b>	<b>0.997 <math>\pm</math> 0.021</b>	<b>0.928 <math>\pm</math> 0.045</b>	<b>0.985 <math>\pm</math> 0.015</b>	<b>0.853 <math>\pm</math> 0.194</b>
434	RankRefine ( $\beta$ )	0.977 $\pm$ 0.024	1.002 $\pm$ 0.022	0.968 $\pm$ 0.016	0.985 $\pm$ 0.020	0.952 $\pm$ 0.055
435	Projection ( $\beta$ )	1.056 $\pm$ 0.171	1.119 $\pm$ 0.100	1.007 $\pm$ 0.026	0.992 $\pm$ 0.039	0.949 $\pm$ 0.082
436	Dataset Name	Half Life	FreeSolv	PPBR	Solubility	VDss
437	Claude4 PRA (%)	52.36 $\pm$ 2.22	72.47 $\pm$ 1.40	51.73 $\pm$ 3.25	60.41 $\pm$ 2.45	60.85 $\pm$ 2.12
438	BAYES-ECR ( $\beta$ )	<b>0.955 <math>\pm</math> 0.058</b>	<b>0.977 <math>\pm</math> 0.011</b>	<b>0.965 <math>\pm</math> 0.042</b>	0.990 $\pm$ 0.010	<b>0.850 <math>\pm</math> 0.191</b>
439	RankRefine ( $\beta$ )	0.981 $\pm$ 0.010	1.070 $\pm$ 0.052	0.979 $\pm$ 0.019	<b>0.988 <math>\pm</math> 0.016</b>	0.951 $\pm$ 0.056
440	Projection ( $\beta$ )	1.124 $\pm$ 0.279	1.153 $\pm$ 0.126	1.031 $\pm$ 0.039	0.995 $\pm$ 0.014	0.959 $\pm$ 0.038

Table 1: Results using LLMs (ChatGPT5, Claude4) as external rankers, measured in  $\beta$  over ten train/test splits with  $N = 50, k = 20$ . PRA stands for pairwise ranking accuracy. Best result for each dataset-ranker is bolded. BAYES-ECR generally outperforms other refinement methods.

Dataset Name	Half Life	FreeSolv	PPBR	Solubility	VDss
$\Delta_y \cdot 10^{-7}$	0.028 $\pm$ 0.028	0.030 $\pm$ 5.71	8.51 $\pm$ 1.42	0.029 $\pm$ 1.33	3.24 $\pm$ 2.96

Table 2: Mean absolute difference between the refined predictions of RankRefine and BAYES-ECR with Gaussianity assumptions. The values are around the range of machine epsilon for a floating point precision, confirming that RankRefine is a special case of BAYES-ECR.

references (50,000 pairs in total) in parallel is  $152.4 \pm 54.7$  seconds. Therefore, BAYES-ECR can be run efficiently on consumer-grade computes, and LLM inference is not a significant bottleneck.

#### 4.6 RANKREFINE IS A SPECIAL CASE OF BAYES-ECR.

We showed that RankRefine (Wijaya et al., 2025) is a special case of BAYES-ECR in Proposition 3.3. We verify this empirically by comparing enhanced predictions from RankRefine and BAYES-ECR under the Gaussianity assumptions. Across five ADMET tasks, the mean absolute difference  $\Delta_y$  is at the level of machine epsilon (Table 2), confirming that Rankrefine is a special case of BAYES-ECR with Gaussianity assumptions.

## 5 LIMITATIONS

- **Noise model for the oracle ranker.** We inject pairwise flips uniformly to reach the desired oracle ranker accuracy. Real rankers (human or LLM) might exhibit systematic biases, which could negatively affect the Bayesian inference at lower ranker accuracies.
- **Dependence on regressor uncertainty.** BAYES-ECR requires a base regressor that can quantify uncertainty, which many off-the-shelf regressors lack. A straightforward solution is to use ensembling or Monte Carlo dropout.
- **Reference-set selection.** Regression improvements may depend on the composition of the reference set (coverage, label diversity, resolution). Future works can explore more sophisticated, sequential sampling strategy to select the references.

## 6 CONCLUSION

We introduced BAYES-ECR, framing post training regression enhancement as a problem of Bayesian inference rather than heuristic fusion. This perspective not only generalizes prior state-of-the-art, but also reveals subtle failure modes, explains when and why performance degrades, and offers principled solutions through temperature calibration and gating. Empirically, BAYES-ECR delivers consistent gains across 12 diverse datasets, operates effectively with both oracle and noisy LLM rankers, and runs under 1 ms per query. We believe BAYES-ECR can serve as a practical bridge between label-scarce applications and the growing availability of pairwise signals from LLMs.

486 REFERENCES  
487

488 Alice EA Allen and Alexandre Tkatchenko. Machine learning of material properties: Predictive and  
489 interpretable multilinear models. *Science advances*, 8(18):eabm7185, 2022.

490 Anthropic. Anthropic api guide. <https://docs.anthropic.com/en/api/overview>,  
491 2025. Accessed: 2025-09-16.

492

493 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method  
494 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

495

496 Manuela Cattelan. Models for paired comparison data: A review with emphasis on dependent data.  
497 *Statistical Science*, pp. 412–433, 2012.

498

499 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
500 reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

501

502 William G Cochran and Sarah Porter Carroll. A sampling investigation of the efficiency of weighting  
503 inversely as the estimated variance. *Biometrics*, 9(4):447–459, 1953.

504

505 Paulo Cortez. Student Performance. UCI Machine Learning Repository, 2014.

506

507 Li Di, Christopher Keefer, Dennis O Scott, Timothy J Strelevitz, George Chang, Yi-An Bi, Yurong  
508 Lai, Jonathon Duckworth, Katherine Fenner, Matthew D Troutman, et al. Mechanistic insights  
509 from comparing intrinsic clearance values between human liver microsomes and hepatocytes to  
guide drug design. *European journal of medicinal chemistry*, 57:441–448, 2012.

510

511 Google. Gemini developer api. <https://ai.google.dev/gemini-api/docs>, 2025. Ac-  
cessed: 2025-09-16.

512

513 Taicheng Guo, Bozhao Nan, Zhenwen Liang, Zhichun Guo, Nitesh Chawla, Olaf Wiest, Xiangliang  
514 Zhang, et al. What can large language models do in chemistry? a comprehensive benchmark on  
515 eight tasks. *Advances in Neural Information Processing Systems*, 36:59662–59688, 2023.

516

517 Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document  
analysis and recognition*, volume 1, pp. 278–282. IEEE, 1995.

518

519 Eva JI Hoeijmakers, Bibi Martens, Babs MF Hendriks, Casper Mihl, Razvan L Miclea, Walter H  
520 Backes, Joachim E Wildberger, Frank M Zijta, Hester A Gietema, Patricia J Nelemans, et al.  
521 How subjective ct image quality assessment becomes surprisingly reliable: pairwise comparisons  
522 instead of likert scale. *European Radiology*, 34(7):4494–4503, 2024.

523

524 Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Connor Co-  
525 ley, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning  
526 datasets and tasks for drug discovery and development. *Advances in neural information process-  
527 ing systems*, 2021.

528

529 Qiaohao Liang, Aldair E Gongora, Zekun Ren, Armi Tiihonen, Zhe Liu, Shijing Sun, James R De-  
530 neault, Daniil Bash, Flore Mekki-Berrada, Saif A Khan, et al. Benchmarking the performance  
531 of bayesian optimization across multiple experimental materials science domains. *npj Computational  
Materials*, 7(1):188, 2021.

532

533 Franco Lombardo and Yankang Jing. In silico prediction of volume of distribution in humans.  
534 extensive data set and the exploration of linear and nonlinear methods coupled with molecular  
535 interaction fields descriptors. *Journal of chemical information and modeling*, 56(10):2042–2052,  
2016.

536

537 Alexander Ly, Maarten Marsman, Josine Verhagen, Raoul PPP Grasman, and Eric-Jan Wagenmak-  
538 ers. A tutorial on fisher information. *Journal of Mathematical Psychology*, 80:40–55, 2017.

539

David L Mobley and J Peter Guthrie. Freesolv: a database of experimental and calculated hydration  
free energies, with input files. *Journal of computer-aided molecular design*, 28:711–720, 2014.

540 Frederick Mosteller. Remarks on the method of paired comparisons: III. a test of significance for  
 541 paired comparisons when equal standard deviations and equal correlations are assumed. *Psychometrika*, 16(2):207–218, 1951.

542

543 Kevin P Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.

544

545 R Scott Obach, Franco Lombardo, and Nigel J Waters. Trend analysis of a database of intravenous  
 546 pharmacokinetic parameters in humans for 670 drug compounds. *Drug Metabolism and Disposi-*  
 547 *tion*, 36(7):1385–1405, 2008.

548

549 OpenAI. Openai api reference. [https://platform.openai.com/docs/](https://platform.openai.com/docs/api-reference/introduction)  
 550 api-reference/introduction, 2025. Accessed: 2025-09-16.

551

552 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
 553 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
 554 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
 27730–27744, 2022.

555

556 Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Le Yan, Jiaming Shen, Tianqi Liu,  
 557 Jialu Liu, Donald Metzler, et al. Large language models are effective text rankers with pairwise  
 558 ranking prompting. In *Findings of the Association for Computational Linguistics: NAACL 2024*,  
 559 pp. 1504–1518, 2024.

560

561 Jennifer Routh, Sharmini Julita Paramasivam, Peter Cockcroft, Sarah Wood, John Remnant,  
 562 Cornélie Westermann, Alison Reid, Patricia Pawson, Sheena Warman, Vishna Devi Nadarajah,  
 563 et al. Rating and ranking preparedness characteristics important for veterinary workplace clin-  
 564 ical training: a novel application of pairwise comparisons and the elo algorithm. *Frontiers in  
 Medicine*, 10:1128058, 2023.

565

566 David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-  
 567 propagating errors. *nature*, 323(6088):533–536, 1986.

568

569 Benjamin Sanchez-Lengeling and Alán Aspuru-Guzik. Inverse molecular design using machine  
 570 learning: Generative models for matter engineering. *Science*, 361(6400):360–365, 2018.

571

572 Adil Shamim. Cost of international education, 2025. URL <https://www.kaggle.com/datasets/adilshamim8/cost-of-international-education>.

573

574 Murat Cihan Sorkun, Abhishek Khetan, and Süleyman Er. Aqsoldb, a curated reference set of  
 575 aqueous solubility and 2d descriptors for a diverse set of compounds. *Scientific data*, 6(1):143,  
 2019.

576

577 Atharva Soundankar. Smart farming sensor data for yield prediction, 2025.  
 578 URL <https://www.kaggle.com/datasets/atharvasoundankar/smart-farming-sensor-data-for-yield-prediction>.

579

580 Michael Sun, Gang Liu, Weize Yuan, Wojciech Matusik, and Jie Chen. Foundation molecular  
 581 grammar: Multi-modal foundation models induce interpretable molecular graph languages. In  
 582 *International Conference on Machine Learning*. PMLR, 2025.

583

584 LL Thurstone. A law of comparative judgment. *Psychological Review*, 34(4):273–286, 1927.

585

586 Ning-Ning Wang, Jie Dong, Yin-Hua Deng, Min-Feng Zhu, Ming Wen, Zhi-Jiang Yao, Ai-Ping Lu,  
 587 Jian-Bing Wang, and Dong-Sheng Cao. Adme properties evaluation in drug discovery: prediction  
 588 of caco-2 cell permeability using a combination of nsga-ii and boosting. *Journal of chemical  
 information and modeling*, 56(4):763–773, 2016.

589

590 David Weininger. Smiles, a chemical language and information system. 1. introduction to method-  
 591 ology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36,  
 1988.

592

593 Kevin Tirta Wijaya, Minghao Guo, Michael Sun, Hans-Peter Seidel, Wojciech Matusik, and Vahid  
 Babaei. Two-stage pretraining for molecular property prediction in the wild. *arXiv preprint  
 arXiv:2411.03537*, 2024.

594 Kevin Tirta Wijaya, Michael Sun, Minghao Guo, Hans-Peter Seidel, Wojciech Matusik, and Vahid  
 595 Babaei. Post hoc regression refinement via pairwise rankings. *arXiv preprint arXiv:2508.16495*,  
 596 2025.

597 Tong Wu, Guandao Yang, Zhibing Li, Kai Zhang, Ziwei Liu, Leonidas Guibas, Dahua Lin, and  
 598 Gordon Wetzstein. Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. In  
 599 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22227–  
 600 22238, 2024.

601 Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S  
 602 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learn-  
 603 ing. *Chemical science*, 9(2):513–530, 2018.

604 Le Yan, Zhen Qin, Honglei Zhuang, Rolf Jagerman, Xuanhui Wang, Michael Bendersky, and Harrie  
 605 Oosterhuis. Consolidating ranking and relevance predictions of large language models through  
 606 post-processing. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Lan-  
 607 guage Processing*, pp. 410–423, 2024.

608 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
 609 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and  
 610 chatbot arena. *Advances in neural information processing systems*, 36:46595–46623, 2023.

611 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
 612 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv  
 613 preprint arXiv:1909.08593*, 2019.

614  
 615  
 616  
 617  
 618  
 619  
 620  
 621  
 622  
 623  
 624  
 625  
 626  
 627  
 628  
 629  
 630  
 631  
 632  
 633  
 634  
 635  
 636  
 637  
 638  
 639  
 640  
 641  
 642  
 643  
 644  
 645  
 646  
 647

648 **A APPENDIX**649 **A.1 DETAILED PROOFS**

650 **Lemma 3.1.** *Proof.* The regressor log-likelihood contributes  $-\frac{1}{2\sigma_{\text{re}}^2}(\hat{y}_0^{\text{re}} - y_0)^2$  with curvature  
 651  $-1/\sigma_{\text{re}}^2 < 0$ . Each BT term adds  $\log s(y_0 - y_i)$  or  $\log 1 - s(y_0 - y_i)$ , whose second derivative is  
 652  $-s(1 - s) \leq 0$ . Sum of concave terms and a concave prior is concave; the Gaussian term's  
 653 strictly negative curvature makes the sum strictly concave.

654 **Lemma 3.4.** *Proof.*  $\mathcal{L}'_{\text{BT}}(y)$  (Proposition 3.3.) is strictly decreasing, therefore,  $\mathcal{L}_{\text{BT}}(y)$  is strictly  
 655 concave. Setting  $\mathcal{L}'_{\text{BT}}(y) = 0$  gives the equation. Since  $F(y) = \sum_i u_i(y)$  is continuous and strictly  
 656 increasing with  $F(-\infty) = 0$ ,  $F(\infty) = k$ , there is a unique solution at the observed count  $m$ .

657 Let  $y_1 \leq \dots \leq y_k$  be the ordered reference set and suppose the true  $y_0 \in (y_m, y_{m+1})$ ,

- 658 • if most  $(y - y_i)$  are in the sigmoid saturated regime,  $F$  behaves like a hard count, so  
 $F(y_0) \approx m$ ,  $F(\tilde{y}_0) = m$ , and  $\hat{y}_0^{\text{ra}} = \tilde{y}_0 \approx y_0$  up to the resolution  $(y_{m+1} - y_m)$ .
- 659 • if many  $(y - y_i)$  lie in the transition region,  $F(y_0) \not\approx m$ ,  $F(\tilde{y}_0) = m$ , and  $\hat{y}_0^{\text{ra}} = \tilde{y}_0 \not\approx y_0$ .  
 This shifts does not induce rank bias as long as  $\tilde{y}_0 \in (y_m, y_{m+1})$ .
- 660 • rank bias arises if the shift exits the interval, i.e.,  $|\tilde{y}_0 - y_0| \geq \min\{y_0 - y_m, y_{m+1} - y_0\}$ ,  
 which is more likely when the number of  $(y - y_i)$  that lie in the transition region is high.

661 **Lemma 3.6.** *Proof.* Let  $I_{\text{reg}} = 1/\sigma_{\text{re}}^2$  be the Fisher information of the regressor term and write the  
 662 total gradient and Fisher information at  $y$  as

$$663 \begin{aligned} g_{\text{tot}}(y) &= g_{\text{reg}}(y) + g_{\text{rank}}(y) \\ 664 &= -\frac{1}{\sigma_{\text{re}}^2}(y - \hat{y}_0^{\text{re}}) + \sum_{i=1}^k (r_i - u_i(y)), \end{aligned} \tag{20}$$

$$665 I_{\text{tot}}(y) = I_{\text{reg}} + I_{\text{rank}}(y). \tag{21}$$

666 Define the local rank target for a single Newton step

$$667 \begin{aligned} \tilde{y}^{\text{ra}}(y) &= y + \frac{g_{\text{rank}}(y)}{I_{\text{rank}}(y)} \\ 668 &= y + \frac{\sum_{i=1}^k (r_i - u_i(y))}{I_{\text{rank}}(y)}. \end{aligned} \tag{22}$$

669 By substituting  $g_{\text{reg}}(y) = I_{\text{reg}} \cdot (\hat{y}_0^{\text{re}} - y)$  and  $g_{\text{rank}}(y) = I_{\text{rank}}(y) \cdot (\tilde{y}^{\text{ra}}(y) - y)$ , a single Newton  
 670 step is an information-weighted average as written in Equation 12.

671 **Lemma 3.7.** *Proof.* Using  $\mathbb{E}[r_i] = (1 - a) + (2a - 1)\mathbb{1}(y_0 > y_i)$  to model flip errors in a noisy  
 672 binary ranker gives

$$673 \mathbb{E}[g_{\text{rank}}(y)] = \sum_{i=1}^k \mathbb{E}[r_i] - \sum_{i=1}^k u_i(y) = k((1 - a) + (2a - 1)p^* - \hat{p}(y)), \tag{23}$$

674 If  $p^* \geq 1/2 \geq \hat{p}(y)$  or  $p^* \leq 1/2 \leq \hat{p}(y)$ , then

$$675 \left| \mathbb{E}[g_{\text{rank}}(y)] \right| \geq k(2a - 1)|p^* - \hat{p}(y)|. \tag{24}$$

676 Substituting  $g_{\text{rank}}(y)$  with  $-1/\sigma_{\text{re}}^2(y - \hat{y}_0^{\text{re}})$ , the expected Newton step is rank-dominated whenever

$$677 (2a - 1)|p^* - \hat{p}(y)| > \frac{|y - \hat{y}_0^{\text{re}}|}{k\sigma_{\text{re}}^2}. \tag{25}$$

702 A.2 USE OF LARGE LANGUAGE MODELS  
703704 We use large language models (LLMs) to (i) polish writing (restructuring sentences and proofread  
705 grammars and typos), (ii) search for related works, (iii) predict pairwise rankings of molecule pairs.  
706 For writing and search, we use ChatGPT5. For pairwise ranking predictions, we use ChatGPT5 and  
707 Claude Sonnet 4 (Copilot version). The prompt that we used for pairwise ranking prediction is as  
708 follow,

709

710 You are an expert molecular reasoning model tasked with  
711 predicting pairwise rankings of molecules based on a described  
712 molecular property of interest (e.g., solubility, polarity, etc.).

713

714 You will be given json files containing the test molecules to be  
715 compared with the reference molecules. The property of interest  
716 is described within the json files with the key "description".

717

\* Your Task - Follow These Steps:

1. Read the dataset's description to understand the molecular  
property being ranked (e.g., "higher solubility", "lower toxicity").  
- Be careful with the measurement unit. For example, a higher number  
in IC50 could mean lower toxicity.2. For each molecule pair (test molecule vs. reference molecule):  
- Use your internal knowledge to infer which molecule ranks higher  
for the property.- You may use structural patterns, substrings, atom types,  
SMARTS-like features, token-level patterns, or other insights  
you may have.- You are encouraged to develop your own heuristics or scoring logic  
using Python.

3. Assign a "pairwise\_rank" to each pair:

1 → test molecule ranks higher than reference, meaning  
test\_property > reference\_property

0 → test molecule ranks lower or equal to reference

Caution! Only care for the value of the property. For example,  
if toxicity is measured in IC50, and test\_molecule IC50 value  
is greater than reference\_molecule IC50 value, you should output 1.

4. Save your results as a new JSON file with similar structure.

736

\* You Are Allowed To:

1. Write your own Python logic.  
2. Use basic Python and string-based pattern recognition  
3. Think step-by-step to develop useful ranking heuristics  
4. Use helper functions from Cheminformatics libraries, as long  
as you do not use them to directly predict the property of interest  
values.

744

\* You Are NOT Allowed To:

1. Use the internet to search for the property values  
2. Use cheminformatics libraries like RDKit to directly predict  
the property of interest values, e.g., solubility of molecule A  
for solubility datasets.  
3. Access files not specified in this prompt

751

\* Each entry in your output JSON should look like this:

"test\_molecule": "smiles": "...",  
"reference\_molecules": [  
 { "id": "...",  
 "smiles": "..."  
 "pairwise\_ranks": 1,

```

756 },
757 that is, put the pairwise_rank predictions inside the
758 reference_molecule. Your output json file should be named
759 "pairwise_ranking_predictions_{split_id}.json"
760
761 * Tips for Better Performance:
762 1. Think aloud: before you begin ranking, describe what is the
763 property of interest and why one molecule might rank higher
764 based on your knowledge
765 2. Use token or substring patterns (e.g., "more OH groups" or
766 "more aromatic rings")
767 3. Define scoring rules: e.g., "count('O') - count('N')"
768 to estimate polarity
769
770
```

### A.3 EFFECTS OF NUMBER OF SAMPLES FOR TEMPERATURE CALIBRATION

In Section 3.3, the temperature  $\tau$  is set to  $\hat{\tau}_{\text{cal}}$ , the calibrated value estimated from the labeled reference set. A natural question is the robustness of this procedure in data-scarce regimes, where only a few labeled references are available. Figure 5 shows  $\hat{\tau}_{\text{cal}}$  as a function of the calibration set size  $k_{\text{cal}}$ . Across the 9 TDC ADMET datasets,  $\hat{\tau}_{\text{cal}}$  converges at around  $k_{\text{cal}} \approx 10$ , indicating that optimal calibration can be obtained with only 10 labeled samples. Furthermore, except for PPBR\_AZ, the estimates obtained with  $k_{\text{cal}} < 10$  are already close to their converged values. This suggests that  $\hat{\tau}_{\text{cal}}$  remains reliable even when very limited calibration data are available. Consistent with this observation, in the experiment where we fix the reference set size at  $k = 50$  but vary  $k_{\text{cal}} \in [3, 50]$  (Figure 6), the resulting performance curves are nearly identical. Together, these results demonstrate that the number of references used for temperature calibration has a negligible effect on overall regression enhancement performance.

### A.4 REGRESSION ENHANCEMENT FOR MULTIPLE TARGETS

The discussion and experiments in Sections 3 and 4 focus on scalar regression tasks. In many practical settings, however, the goal is to predict multiple properties simultaneously, i.e., to produce vector-valued outputs. Extending BAYES-ECR to this setting is straightforward: one can assume independence across output dimensions and apply BAYES-ECR separately to each component of a target vector  $\mathbf{y} \in \mathbb{R}^d$ , yielding  $d$  independent regression enhancement processes.

In practice, the components of  $\mathbf{y}$  are often correlated. For example, in molecular chemistry, aqueous solubility and membrane permeability are often related, as are clearance and half-life. In such cases, dependencies among outputs can be captured by employing a low-rank approximation of the covariance structure, which enables a computationally efficient implementation. Exploring this extension is an interesting direction for future work.

### A.5 MORE RESULTS ON TDC ADMET

We show results for more  $k$  values in Figure 7-Figure 15. In general, BAYES-ECR can improve the regression error (i.e.,  $\beta < 1$ ) with a reference set size as small as  $k = 3$ .

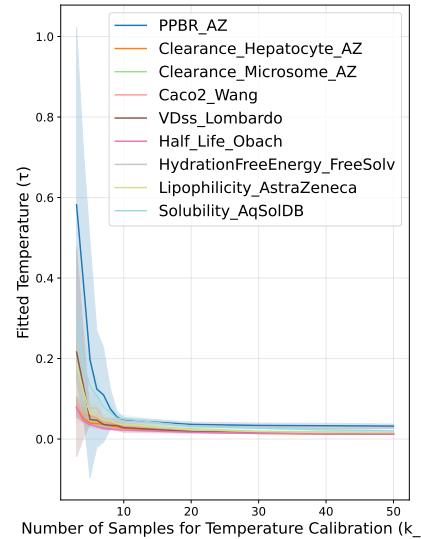


Figure 5: Effects of number of samples used to calibrate the temperature on the value of  $\hat{\tau}_{\text{cal}}$ .

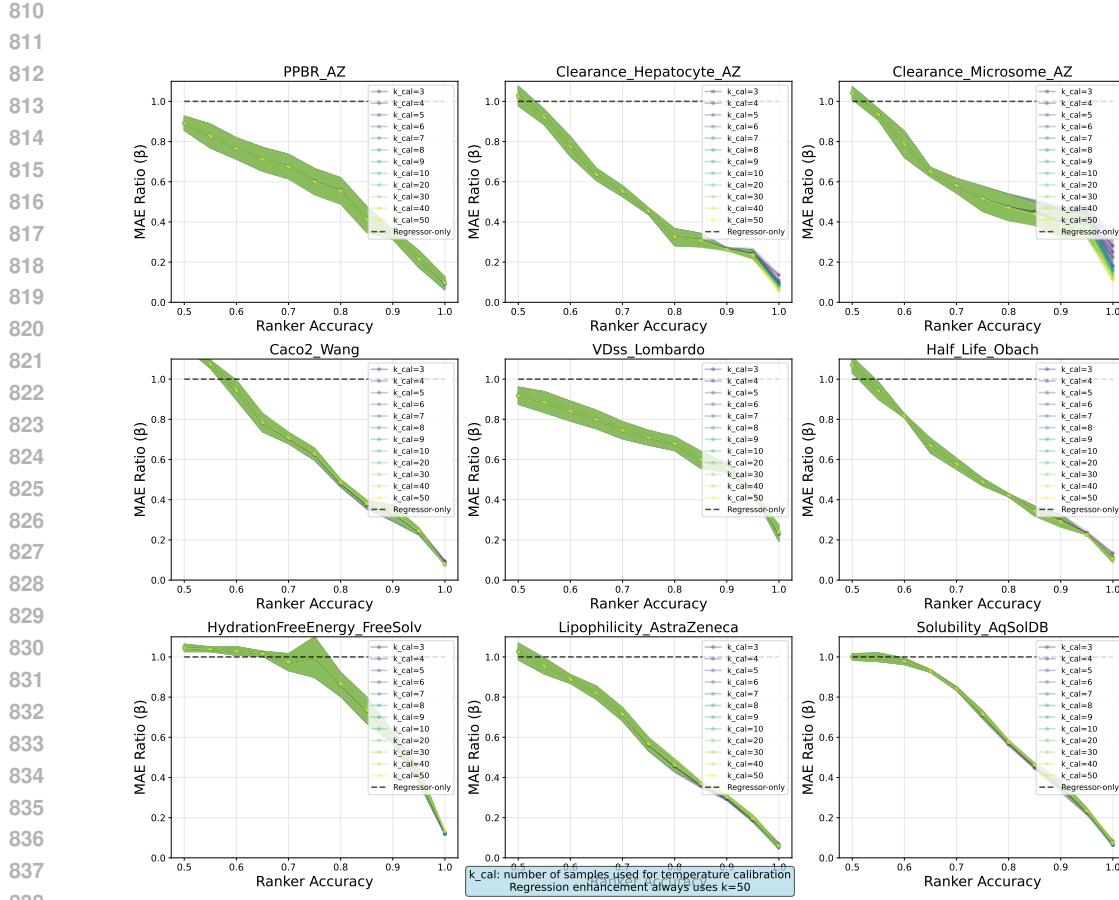


Figure 6: Effects of number of samples on the regression enhancement. We measure  $\beta$  (lower is better) as a function of number of samples used to calibrate the temperature  $k_{cal}$ . The reference set size  $k$  is set to 50.

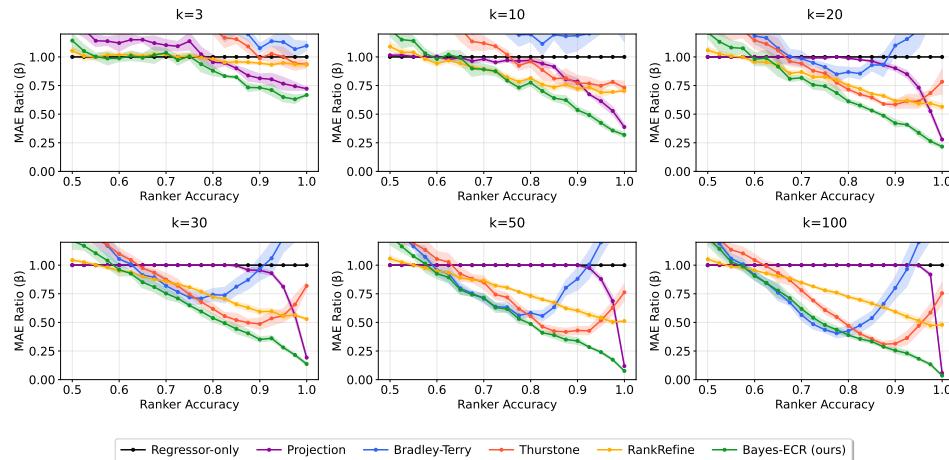


Figure 7: Caco2\_Wang

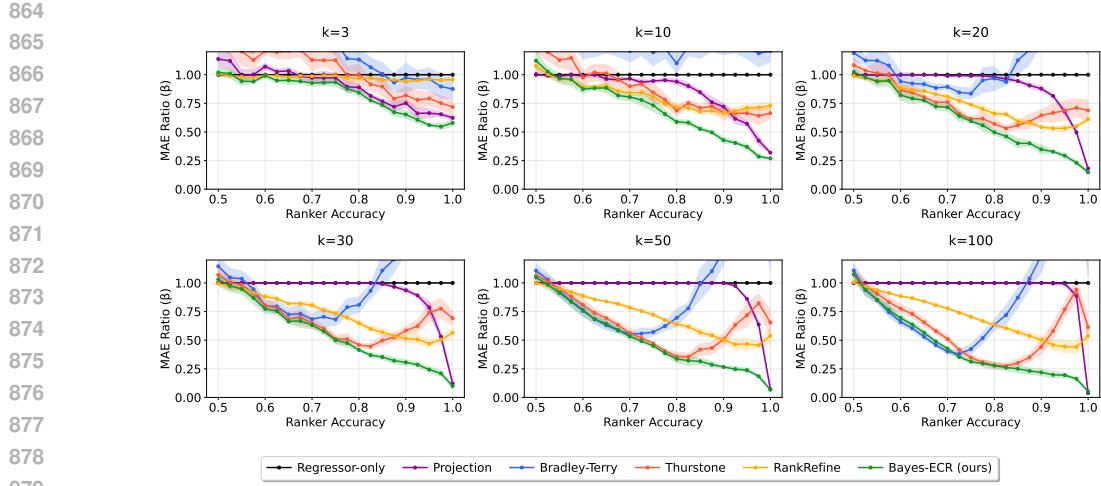


Figure 8: Clearance\_Hepatocyte\_AZ

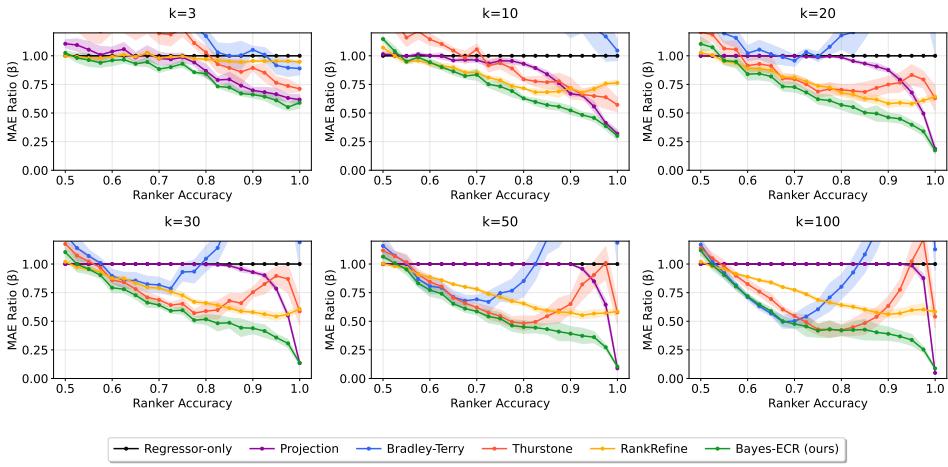


Figure 9: Clearance\_Microsome\_AZ

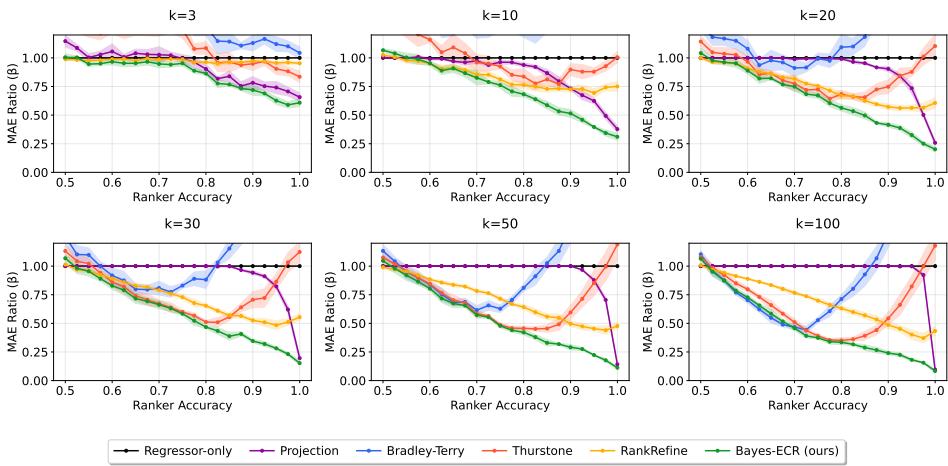


Figure 10: Half\_Life\_Obach

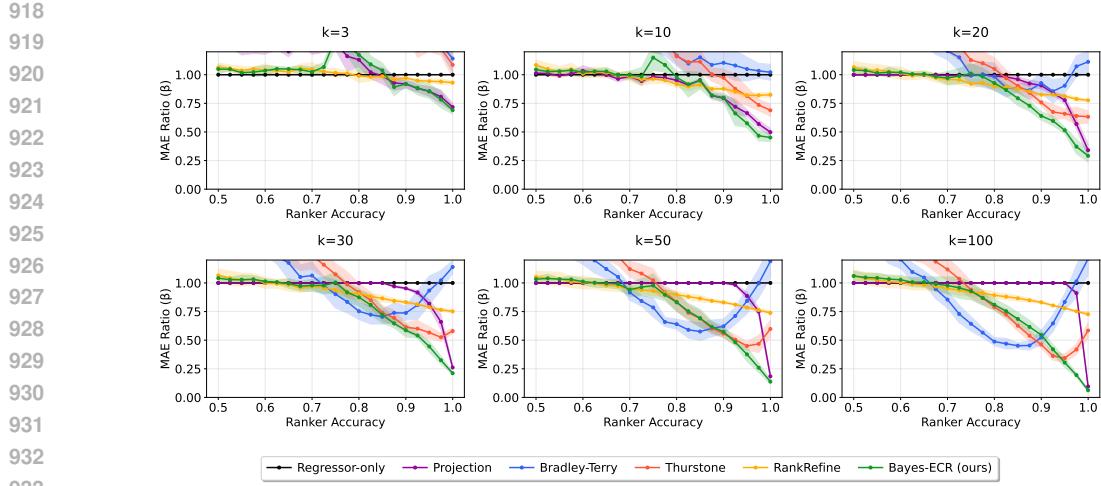


Figure 11: HydrationFreeEnergy\_FreeSolv

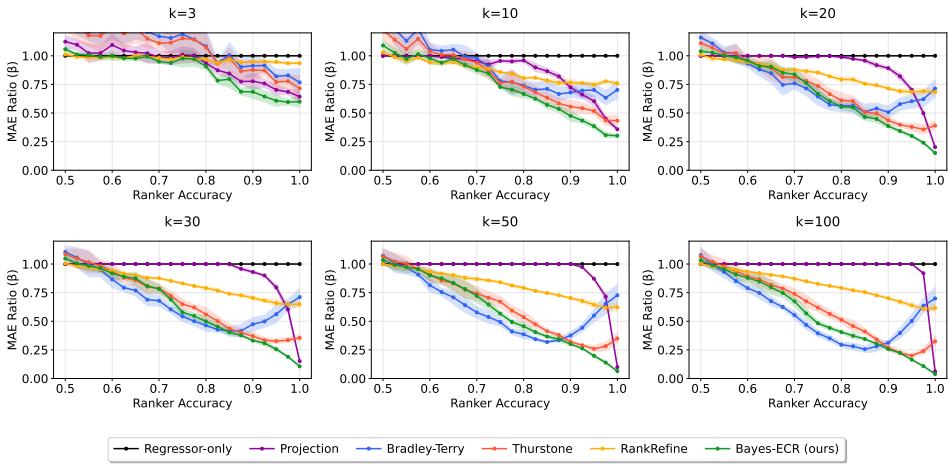


Figure 12: Lipophilicity\_AstraZeneca

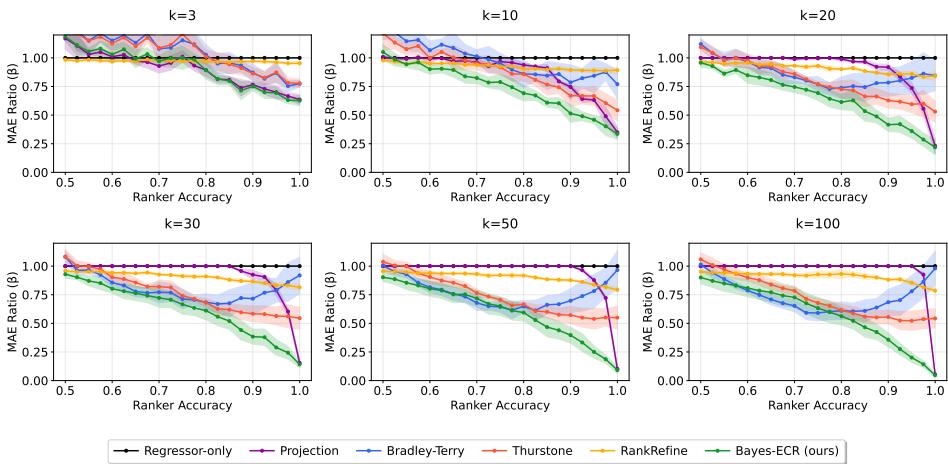


Figure 13: PPBR\_AZ

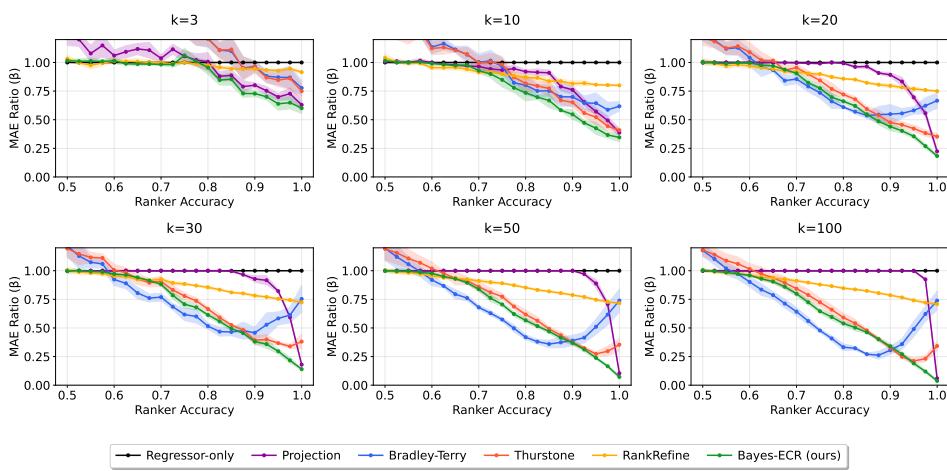


Figure 14: Solubility\_AqSolDB

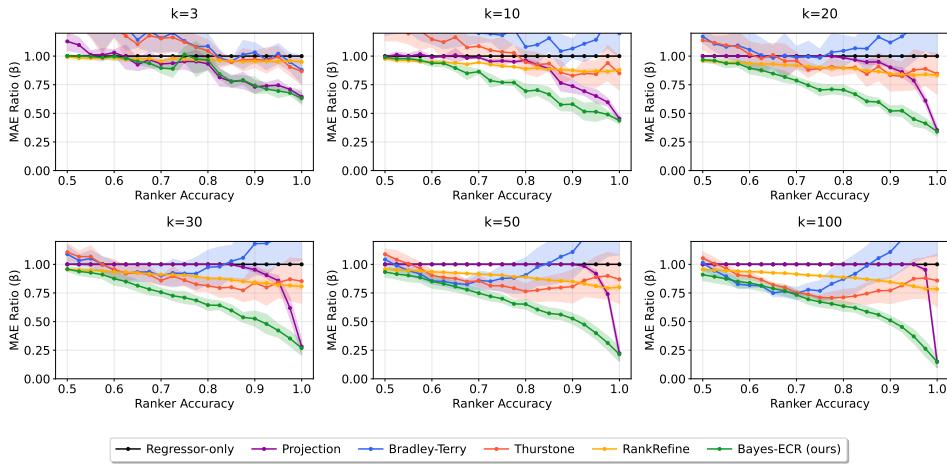


Figure 15: VDss\_Lombardo

1026  
1027

## A.6 REBUTTAL

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

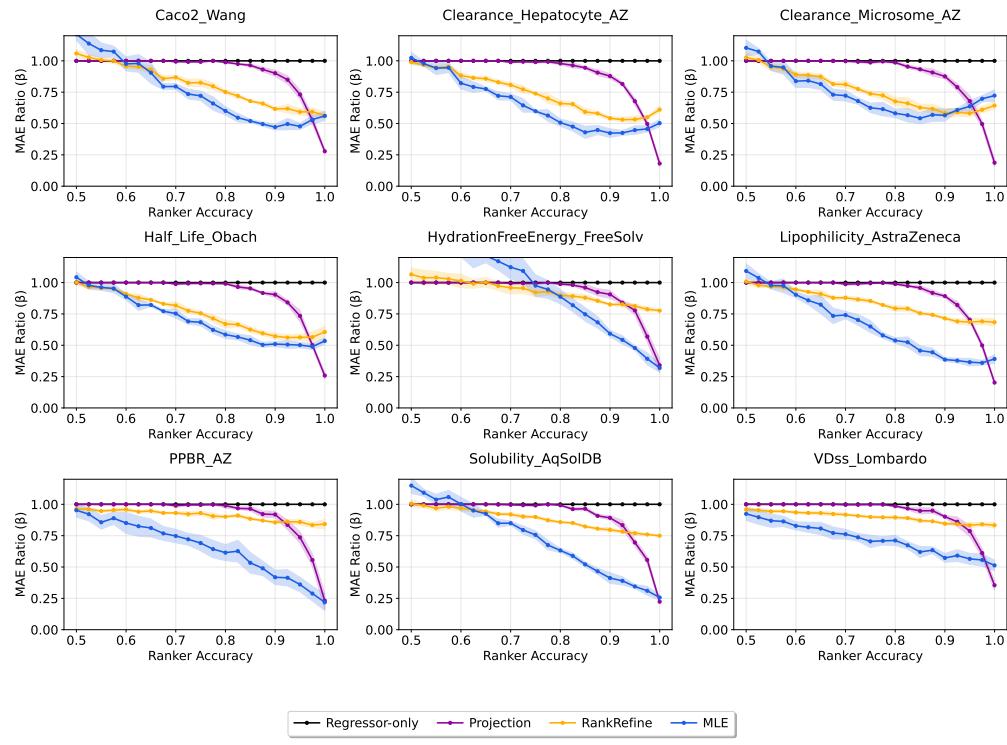


Figure 16: Non-calibrated Bayes-ECR vs. RankRefine and Projection. Reference set size  $k = 20$ . Non-calibrated Bayes-ECR consistently achieves lower error compared to RankRefine and Projection baselines across datasets, often by a large margin

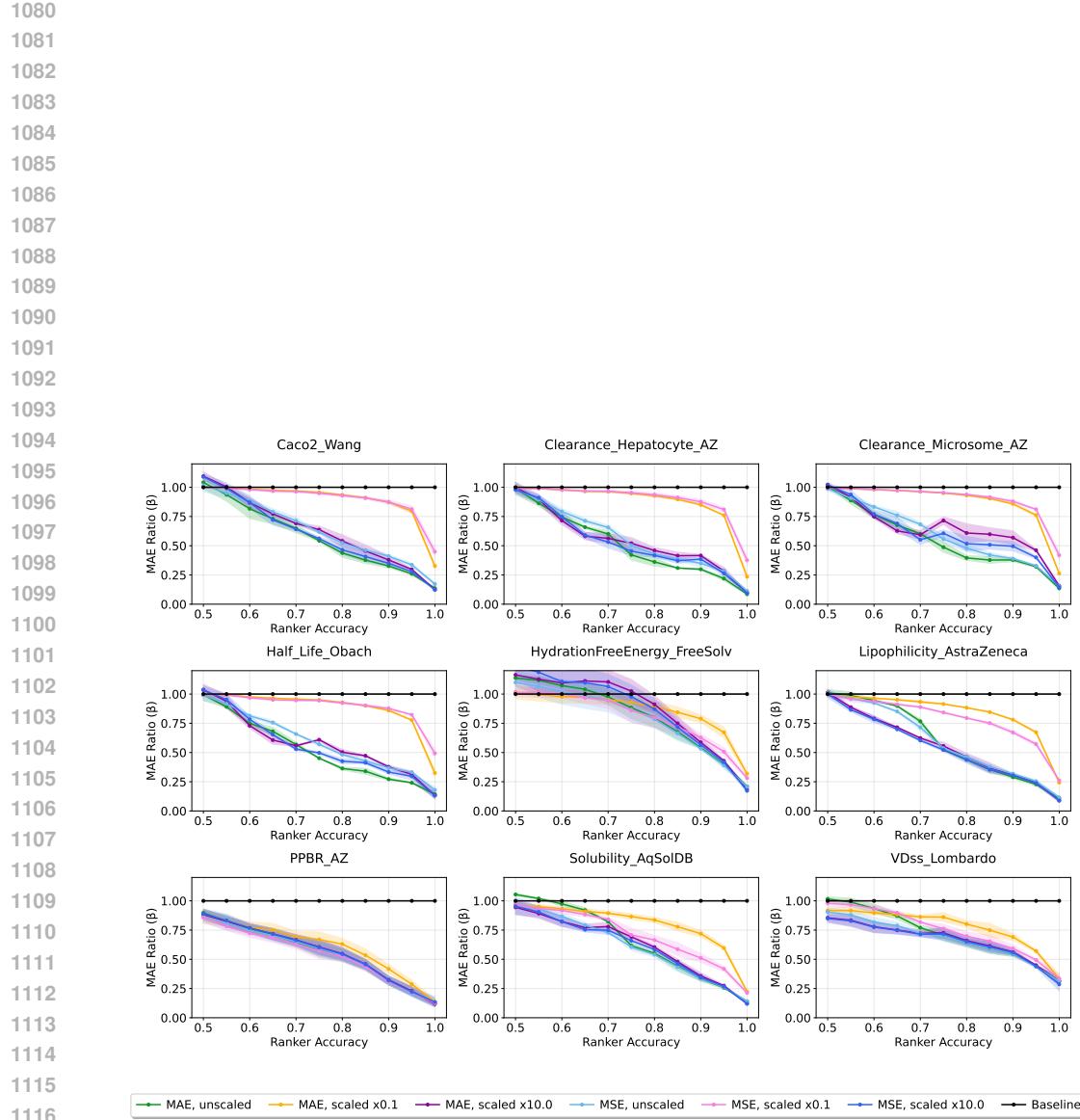


Figure 17: Sensitivity analysis on under or overestimation of regressor uncertainty. Overestimating uncertainty (e.g.,  $s = 10.0$ ) has minimal impact on performance, while underestimating it (e.g.,  $s = 0.1$ ) degrades performance due to excessive confidence in the regressor.

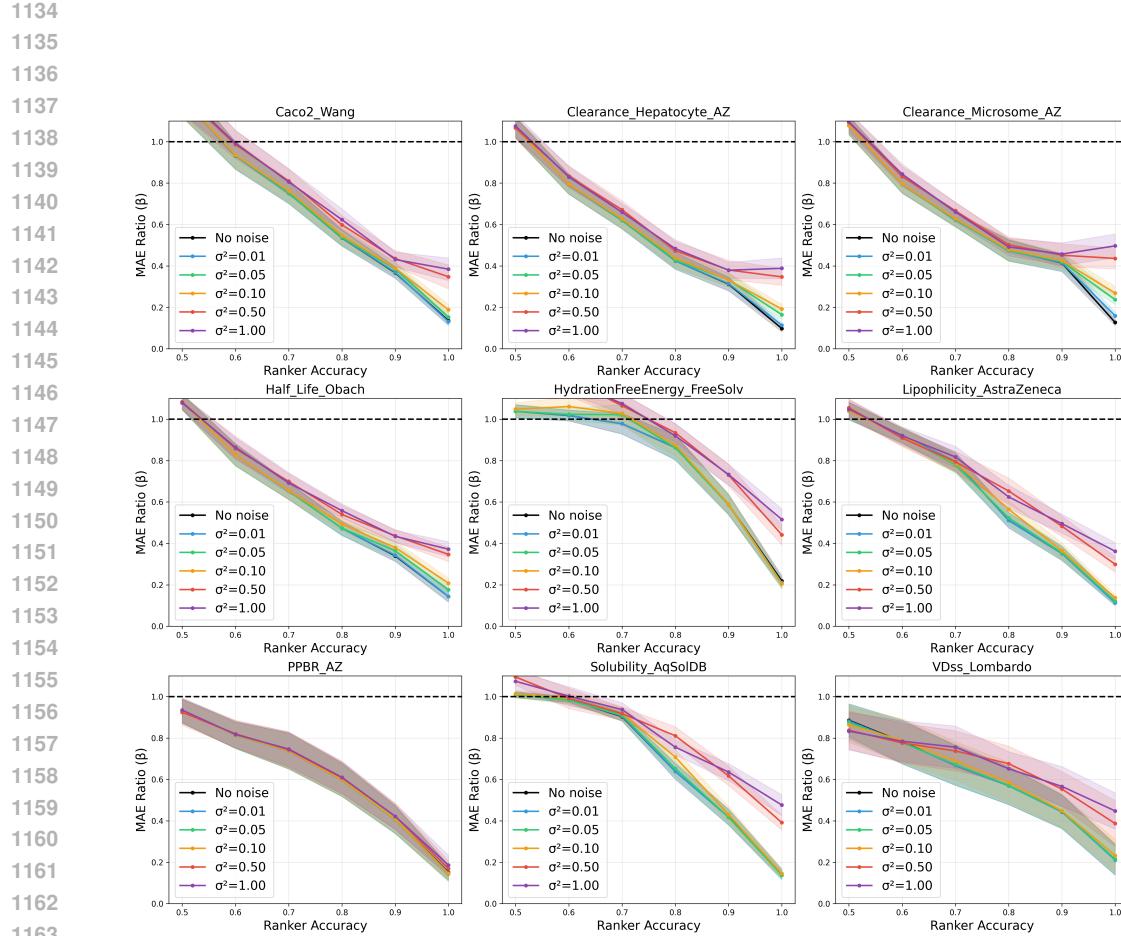


Figure 18: Sensitivity analysis when ranker accuracy is under or overestimated, resulting in non-optimal tau values. Bayes-ECR is robust even when using a suboptimal tau value caused by softgating bias

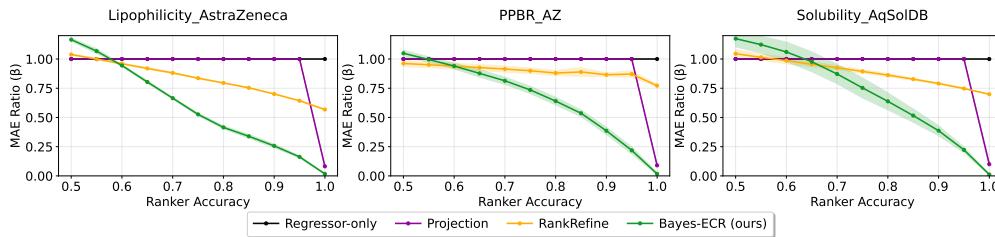
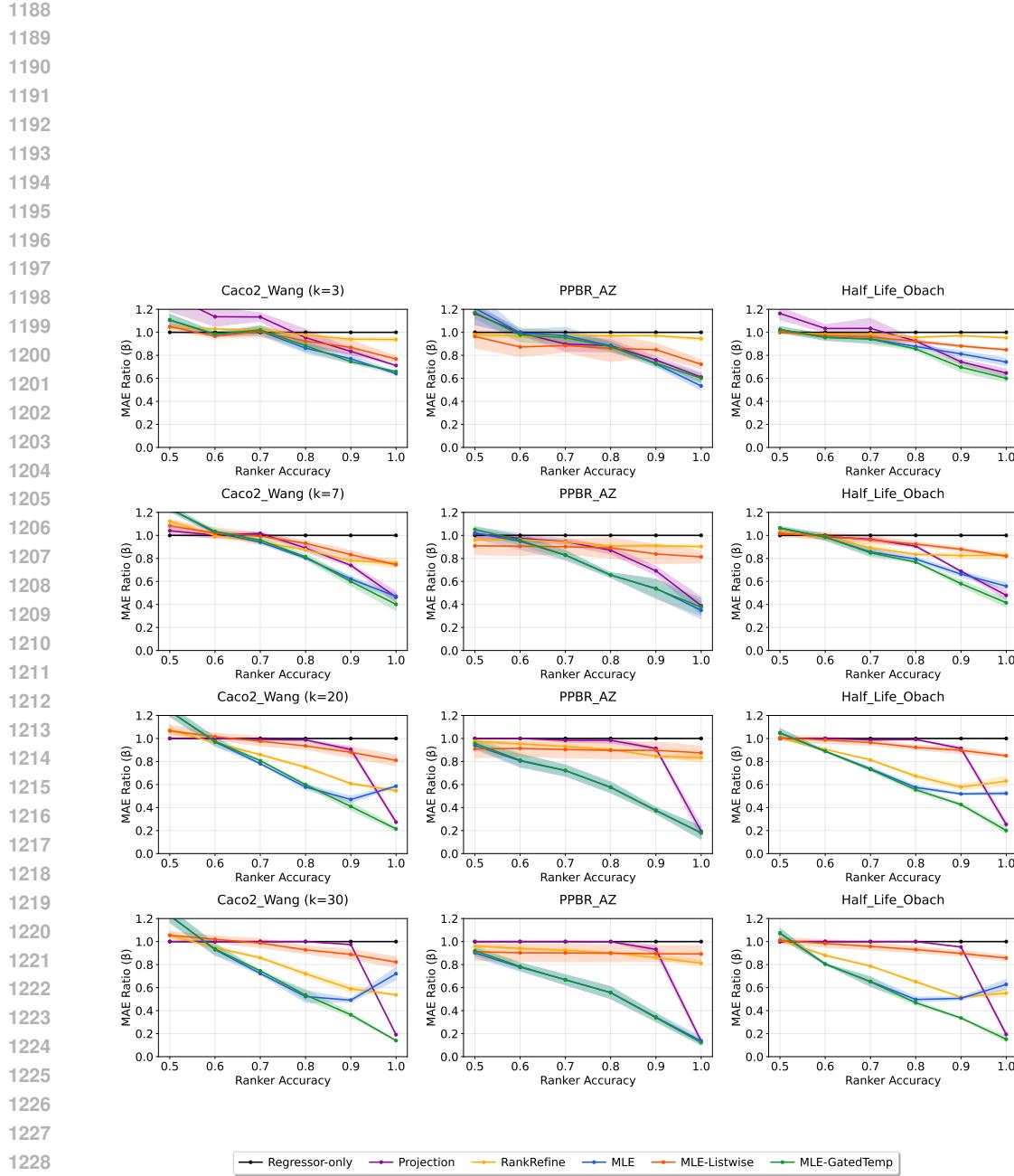


Figure 19: We evaluated our method, Bayes-ECR, on three larger ( $\geq 1000$  training labels) TDC ADME datasets: Lipophilicity, PPBR, and Solubility. With a reference set size of 1000, Bayes-ECR consistently delivered the best overall performance across the entire spectrum of ranker accuracy.



1230 Figure 20: We extend our method, Bayes-ECR, to listwise rankings (MLE-Listwise). For  $k = 3$ ,  
 1231 the listwise variant achieves reasonable performance, confirming that Bayes-ECR can indeed be  
 1232 extended to listwise ranking. However, its performance is not yet on par with our pairwise version.

1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241