

810	A	APPENDIX
811		
812	A.1	GLOSSARY OF ACRONYMS
813	AL	active learning
814	ALP	active learning pipeline
815	AUBC	area under the budget curve
816	DS	dataset
817	ML	machine learning
818	QS	query strategy
819	SOTA	state-of-the-art
820	DNN	deep neural network
821	ETC	extremely randomized trees
822	GBDT	gradient-boosted decision tree
823	k-NN	k-nearest neighbor
824	LR	logistic regression
825	MLP	multi-layer perceptron
826	NB	naïve Bayes
827	PFN	prior-fitted network
828	RF	random forest
829	SVM	support vector machine
830	XGB	XGBoost
831	AAL	adaptive active learning
832	ALBL	active learning by learning
833	BALD	Bayesian active learning by disagreement
834	CER	combined error reduction
835	CLUE	clustering uncertainty-weighted embeddings
836	CluMS	cluster margin
837	CoreSet	CoreSet
838	DWUS	density weighted uncertainty sampling
839	EER	expected error reduction
840	EMC	expected model change
841	ES	entropy sampling
842	EU	epistemic uncertainty sampling
843	EVR	expected variance reduction
844	FALCUN	fast active learning by contrastive uncertainty
845	FIVR	Fisher information variance reduction
846	GRAPH	graph density
847	HIER	hierarchical sampling
848	LC	least-confident sampling
849	k-means	k-means sampling
850	MarginDensity	pre-clustering and margin sampling
851	MaxEnt	maximum entropy
852	MaxER	maximum error reduction
853	MinMS	minimum margin sampling
854	MLI	minimum loss increase
855	MMC	maximum model change
856	MS	margin sampling
857	PowBALD	power-set BALD
858	PowMS	power-set margin sampling
859	QBC	query-by-committee
860	QBC VR	QBC VR
861	QUIRE	querying informative and representative examples
862	Rand	random sampling
863	TypClu	typical clustering
	VR	variance reduction

A.2 COMPARISON TO EXISTING BENCHMARKS FOR TABULAR DATA

In the following, we present an extensive table which compares ALPBench with existing active learning benchmarks. The QS and learning algorithms are ordered by their year of appearance. In Table 3, we additionally present a detailed version of Table 2 in the main paper, which shows which exact QS and learners were implemented in the benchmarks.

Query Strategy	Year	Yang et al. (2018)	Zhan et al. (2021)	Bahri et al. (2022a)	Lu et al. (2023)	ALPBench
ES Shannon (1948)	1948	✓	✓	✓	✓	✓
QBC Seung et al. (1992)	1992	×	✓	×	✓	✓
VR Cohn (1993)	1993	×	✓	×	✓	✓
LC Lewis and Gale (1994)	1994	×	✓	✓	✓	✓
FIVR Zhang (2000)	2000	✓	×	×	×	×
MS Scheffer et al. (2001)	2001	×	✓	✓	✓	✓
EER Roy and McCallum (2001)	2001	✓	✓	×	✓	✓
MaxER Guo and Greiner (2007)	2007	✓	×	×	×	×
CER Guo and Schuurmans (2007b)	2007b	✓	×	×	×	×
EVR Schein and Ungar (2007)	2007	✓	×	×	×	×
EMC Settles et al. (2007)	2007	×	✓	×	×	×
MLI Hoi et al. (2008)	2008	✓	×	×	×	×
BALD Houlisby et al. (2011)	2011	×	×	×	×	×
MMC Cai et al. (2017)	2017	✓	×	×	×	×
MaxEnt Gal et al. (2017)	2017	×	×	✓	×	✓
QBC VR Beluch et al. (2018)	2018	×	×	✓	×	✓
EU Nguyen et al. (2019)	2019	×	×	×	×	✓
PowMS Kirsch et al. (2021)	2021	×	×	✓	×	✓
MinMS Jiang and Gupta (2021)	2021	×	×	✓	×	✓
k-means Kang et al. (2004)	2004	×	✓	×	×	✓
HIER Dasgupta and Hsu (2008)	2008	×	✓	×	✓	×
CoreSet Sener and Savarese (2018)	2018	×	×	✓	✓	✓
TypClu Hacohen et al. (2022)	2022	×	×	✓	×	✓
MarginDensity Nguyen and Smeulders (2004)	2004	×	×	✓	×	×
DWUS Settles and Craven (2008)	2008	×	✓	×	✓	×
QUIRE Huang et al. (2010)	2010	×	✓	×	✓	×
GRAPH Ebert et al. (2012)	2012	×	✓	×	✓	×
AAL Li and Guo (2013)	2013	✓	×	×	×	×
ALBL Hsu and Lin (2015)	2015	×	✓	×	✓	×
CluMS Citovsky et al. (2021)	2021	×	×	✓	×	✓
CLUE Prabhu et al. (2021)	2021	×	×	×	×	✓
FALCUN Gilhuber et al. (2024)	2024	×	×	×	×	✓
Learning Algorithm	Year	Yang et al. (2018)	Zhan et al. (2021)	Bahri et al. (2022a)	Lu et al. (2023)	ALPBench
LR Berkson (1944)	1944	✓	×	×	×	✓
k-NN Fix and Hodges (1952)	1952	×	×	×	×	✓
MLP Werbos (1974)	1974	×	×	✓	×	✓
NB Kononenko (1990)	1990	×	×	×	×	✓
SVM Boser et al. (1992)	1992	×	✓	×	✓	✓
RF Breiman (2001)	2001	×	×	×	×	✓
ETC Geurts et al. (2006)	2006	×	×	×	×	✓
XGB Chen and Guestrin (2016)	2016	×	×	×	×	✓
Catboost Drogush et al. (2018)	2018	×	×	×	×	✓
TabNet Arik and Pfister (2021)	2021	×	×	×	×	✓
TabPFN Hollmann et al. (2023)	2023	×	×	×	×	✓

Table 3: Comparison of the scopes of ALPBench and previous benchmarks for tabular data.

		Yang et al. (2018)	Zhan et al. (2021)	Bahri et al. (2022a)	Lu et al. (2023)	Ours
Qs	Info.	ES, MaxER, MMC, FIVR, EER, CER, EVR, MLI	ES, QBC, VR, LC, MS, EER, EVR	ES, LC, MS, BALD, MaxEnt, QBC VR, PowMS, MinMS, PowBALD	ES, QBC, VR, LC, MS, EER	ES, QBC, VR, LC, MS, EER, BALD, MaxEnt, QBC VR, EU, PowMS, MinMS, PowBALD
	Repr.	-	k-means, HIER	CoreSet, TypClu	HIER, CoreSet	k-means, TypClu, CoreSet
	Hybr.	AAL	DWUS, QUIRE, GRAPH, ALBL	MarginDensity, CluMS	DWUS, QUIRE, GRAPH, ALBL	CluMS, FALCUN, CLUE
Learner	Base	LR	SVM	-	SVM	k-NN, SVM, RF, LR, NB, ETC
	GBDT	-	-	-	-	CatBoost, XGB
	DNN	-	-	MLP	-	MLP, TabNet
	PFN	-	-	-	-	TabPFN
ALP	\sum	9	13	12	12	209
DS	Binary	44	35	35	26	48
	Multi	-	9	34	-	38
	\sum	44	44	69	26	86
AL Setting	1	1	3	1	5	
Metrics	Accuracy	Accuracy	Accuracy	Accuracy	Accuracy, AUC, F1, Prec, Recall, Logloss	

A.3 EXPERIMENTS

In this section, we elaborate in more detail on the experiments that were conducted within our evaluation study.

Datasets. From the 90 datasets from the OpenML-CC18 [Bischl et al. \(2019\)](#) and the TabZilla Benchmark Suite [McElfresh et al. \(2023\)](#) we filtered and excluded the datasets with OpenML IDs 1567, 1169, 41147, and 1493. The first three were filtered out for all settings because they consist of more than 300,000 data points, which would result in a large amount of computing time for the non-info. based Qs. The last dataset with OpenML ID 1493 was filtered out since it consists of 100 classes, which would result in a huge amount of the per iteration budget \mathcal{R} , limiting the number of iterations to a high degree. Further, for the large setting, we wanted to guarantee that at least 10 iterations can be performed until all instances from \mathcal{D}_U are queried. This led to the removal of OpenML IDs 11, 12, 14, 16, 18, 22, 25, 51, 54, 188, 307, 458, 469, 1468, 1501, 40966, and 40979 for this setting. For the preprocessing steps, we proceed as follows. Categorical features are one-hot encoded and missing values are imputed by the mean or mode of the corresponding feature.

Active Learning Setting. As mentioned, we investigate a small and a large setting. Explicitly, the small and large settings are specified by $|\mathcal{D}_L^0| = R = 5 \cdot |\mathcal{Y}|$ and $|\mathcal{D}_L^0| = R = 20 \cdot |\mathcal{Y}|$, respectively, for the given dataset and a total amount of 20 iterations or until all instances from the unlabeled pool \mathcal{D}_U are queried. We choose the factor 5 for the small setting, since then \mathcal{R} matches the one in the (static) small setting in [Bahri et al. \(2022a\)](#). For the large setting, we should have chosen a factor of 100 to be again consistent with [Bahri et al. \(2022a\)](#). However, this seemed unrealistic to us for real-world applications. For some (imbalanced) datasets, it may happen that not every class is at least once represented in \mathcal{D}_L^0 . In these cases, we additionally randomly sample one instance from \mathcal{D}_U per missing class and add them with their corresponding label to \mathcal{D}_L^0 . We run each ALP ten times with different seeds, where the seed defines the $\frac{2}{3}/\frac{1}{3}$ -split of the total dataset \mathcal{D} into $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ as well as the split of $\mathcal{D}_{\text{train}}$ into \mathcal{D}_L and \mathcal{D}_U . Needless to say, the datasets we consider are originally (fully-)labeled datasets. Tailored to the AL setting, we discard the labels for the instances in \mathcal{D}_U and assure that only the oracle \mathcal{O} can access them.

Configuration of Learning Algorithms. In general, we do not perform any hyperparameter optimization (HPO) but rather stick to the default parameters. To contain computational costs, we limit the training time of the learning algorithms. For XGB and Catboost, we reduce the training time by setting the tree method to *hist* and limiting the amount of iterations, respectively. For Catboost and for TabNet, we implement a timeout of three minutes per iteration for the same purpose. This of course may decrease the performance of the learning algorithms and poses a limitation to the generalizability of our empirical study. Further, TabPFN ([Hollmann et al., 2023](#)) can so far only be fitted on a maximum amount of 1,000 instances. Therefore, we uniformly sample 1,000 instances from the current dataset to be fitted on, in case this constraint is violated, similar to [McElfresh et al. \(2023\)](#). For TabPFN and TabNet we modify the implementation for the representation-based and hybrid approaches. Concretely, we extract the output of the encoder from the TabPFN and the activations of the

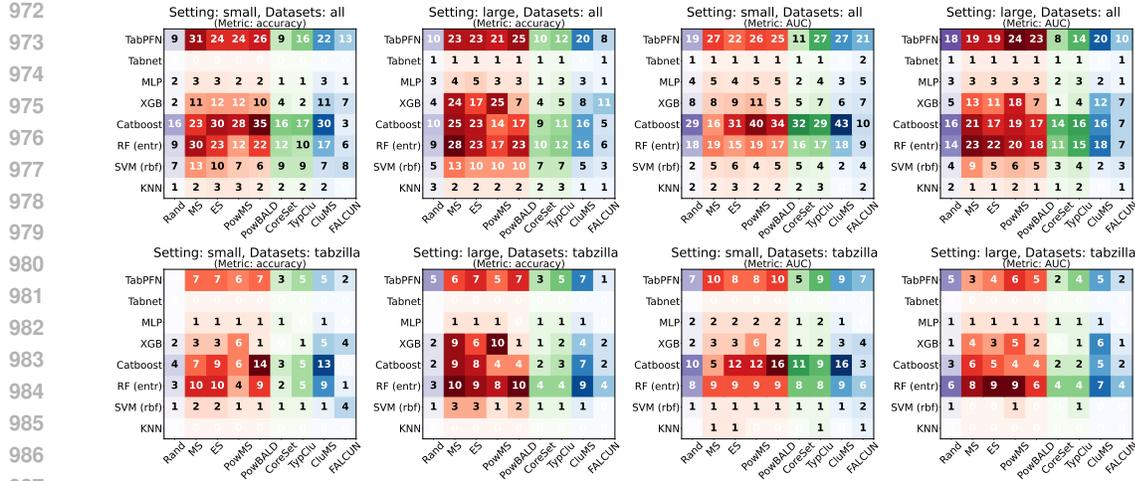


Figure 5: Heatmaps for all ALPs within our evaluation study using AUBC (**accuracy**) as performance measure (first and second column) and AUBC (**AUC**) (third and fourth), separately for **all** (first row) and hybrid Qs are colored in red, green, and blue, respectively, and random sampling is in purple.

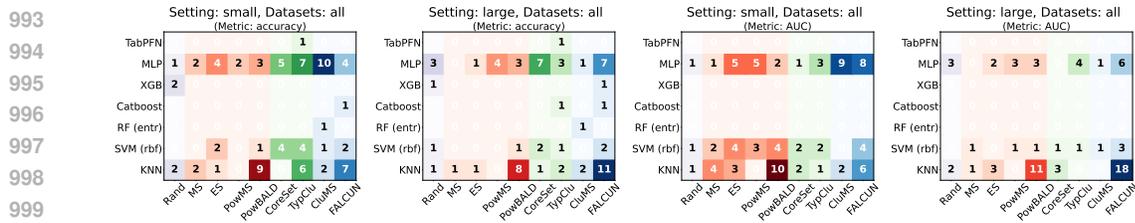


Figure 6: Lose-Heatmaps for all ALPs within our evaluation study using AUBC (**accuracy**) as performance measure (first and second subfigure) and AUBC (**AUC**) (third and fourth) on all datasets **without** statistical significance. The color-coding is consistent with Figure 5.

penultimate layer from TabNet to compute the representativeness of each instance based on its embedding. The exact details can obviously be looked up in our implementation.

Implementation. All experiments were conducted with 2 CPU cores and 8GiB RAM or 16GiB for the small and large settings, respectively, to resemble end-user environments. The HPC nodes for the computations are equipped with two AMD Milan 7763 and 256GiB main memory in total. Runs exceeding these limits have been canceled by the workload manager.

A.4 RESULTS

This section contains more experimental results, comprising more heatmaps and win-matrices distinguishing between binary and multi-class datasets, small and large settings and different metrics. We also present more budget curves for other datasets and learners.

Precisely, we first present heatmaps where we - similar to the main paper - distinguish between small and large settings as well as both metrics AUBC (accuracy) and AUBC (AUC). However, we now compute heatmaps for all datasets (binary and multi-class combined) and for all datasets from the TabZilla Benchmark Suite [McElfresh et al. \(2023\)](#), cf. Figure 5 the first and the second row, respectively. The latter one is a selection of particularly hard or difficult datasets, so we suppose them to be hard for active learning as well.

The main trend of the results of all datasets looks quite similar to the binary datasets in the main paper: Most winning pipelines constitute of TabPFN, Catboost, XGB or RF as learner and information-based QS. However, CluMS is also part of many winning pipelines, especially in the small setting and Rand is quite competitive when considering AUC. For the TabZilla datasets, TabPFN and XGB appear to be not that strong. The QS k-NN and Tabnet (almost) never constitute a winning pipeline and CluMS again is competitive regarding both metrics, especially in the small setting.

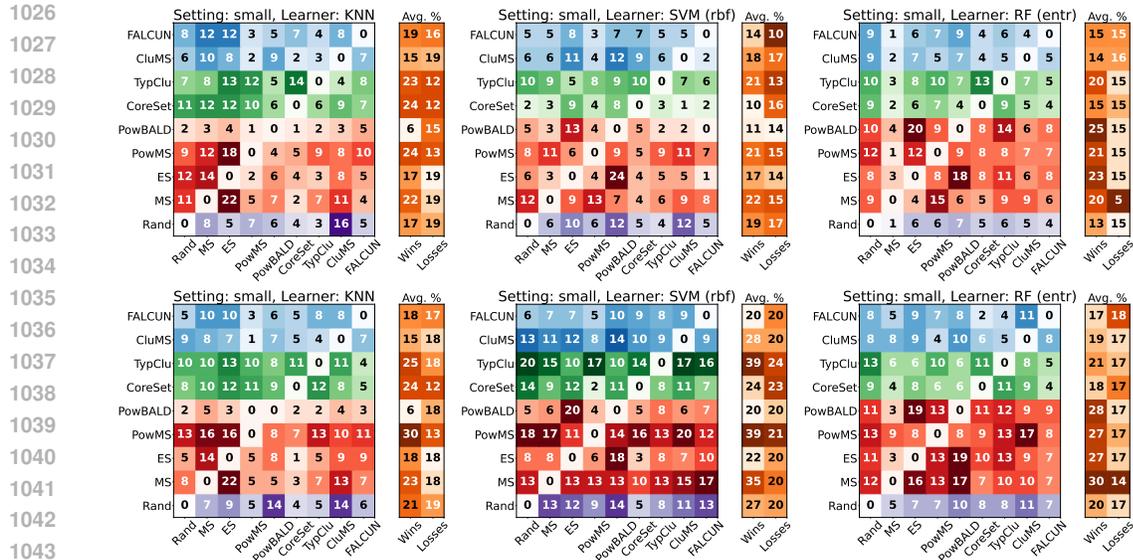


Figure 7: Win-Matrices for k-NN, SVM and RF for the **small** setting on **multi-class** datasets using AUBC (**accuracy**) as performance measure (first row) and AUBC (**AUC**) (second row).

To investigate which ALPs perform particularly poorly, we present *Lose-Heatmaps* in Figure 6, where the losing pipeline replaces ALP_d . Hereby, we do not separate between binary and multi-class datasets and further exclude TabNet as it did not perform at all in our investigated setting. We neglect statistical significance, which may seem an unusual perspective, but it helps to reveal insights into which ALPs exhibit the lowest performance for each dataset. In this figure, we find that it is more important to choose a strong learner than selecting a suitable QS. Concretely, one should avoid MLP or k-NN, and ALPs combining k-NN with PowBALD or MLP with FALCUN or CluMS proved disadvantageous. It might happen, that your learner is not strong, because you maybe want to use a very simple, interpretable model or the data is extremely difficult to learn. In this case it might not be a good idea to rely on any probabilistic estimates but rather choose Rand, as it rarely constitutes to losing pipelines for k-NN and MLP.

In Figure 7 we present win-matrices for the learners k-NN, SVM and RF considering the small setting and evaluating on multi-class datasets. Hereby, we distinguish again between the metrics AUBC (accuracy) and AUBC (AUC). If the metric is chosen as accuracy, we make the following observations. For the k-NN the representation-based and hybrid approaches are very competitive with the information-based strategies. This effect decreases, when SVM is chosen and for the RF the information-based strategies are dominant with MS being extremely robust. In contrast to the RF, Rand is not a too bad choice for k-NN and SVM. Regarding the AUC, TypClu is quite strong for the SVM. For the RF, the information-based strategies are outperforming other QS and in particular MS is strong. Again, we see that the performance of all QSs depend on the chosen learning algorithm.

Further, we present budget curves comparing a subset of 5 different QS for enhanced visual clarity. Precisely, we chose Rand, two representatives for the information-based strategies (MS and power-set BALD (PowBALD)), and one representative for each remaining group, namely CoreSet (CoreSet) and CluMS.

For the large setting, we present budget curves for the datasets with OpenML ID 3 and 1043 in Figure 8. For both datasets, MS is a strong competitor, however CluMS seems to be very strong in the first few iterations. Rand is outperformed by all other strategies, except for the XGB on the first dataset. If the learner achieves high accuracy (as XGB and Catboost do), its probability estimates seem to be reliable and hence information-based strategies are very strong. For the dataset with ID 1043, we observe that CoreSet is initially also quite competitive. If initially the learner has not yet learned too much about the data distribution and achieves also not too good test performance (less than 0.8 accuracy), it might be advantageous to sample representative instances.

In Figure 9, we present budget curves for the datasets with OpenML ID 11 and 51, which both are included in the TabZilla benchmark suite. For the first dataset, one can see that the budget curves for the strong learners RF and TabPFN look quite smooth, especially for TabPFN and also achieve quite high accuracy. The simpler learners k-NN and MLP are struggling more and k-NN even drops in performance in the second half of the active learning procedure. The suitability of different QS again, is quite dependent on the learner: Whereas for the MLP and TabPFN the information-based strategies MS and PowBALD are outperforming the rest, they are the worst when considering k-NN and RF as learners. Regarding the dataset with ID 51, all learners have a hard

1080 time learning the data distribution, as the budget curve is very noisy and also the increases in accuracy are very
1081 marginal, except for the MLP. One can deduce, that this dataset definitely is hard for active learning.

1082 In Figure 10, we consider the small setting and present budget curves for the dataset with OpenML ID 334.
1083 Overall, the budget curves are much less unstable, compared to the large setting. This is expected, as we
1084 start with a very small initial labeled dataset, which makes it really hard to learn the data distribution. The
1085 performance of the different QS differs quite a lot for different learners. CoreSet is very strong if the learning
1086 algorithms is chosen to be k-NN or TabPFN, whereas for both other learners, the information-based strategies
1087 are quite strong. The pipelines consisting of TabPFN as learning algorithm achieve all a much higher accuracy
1088 than the pipelines constituted of the other learners. This highlights the importance of choosing an appropriate
1089 learning algorithm for the given dataset.

1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

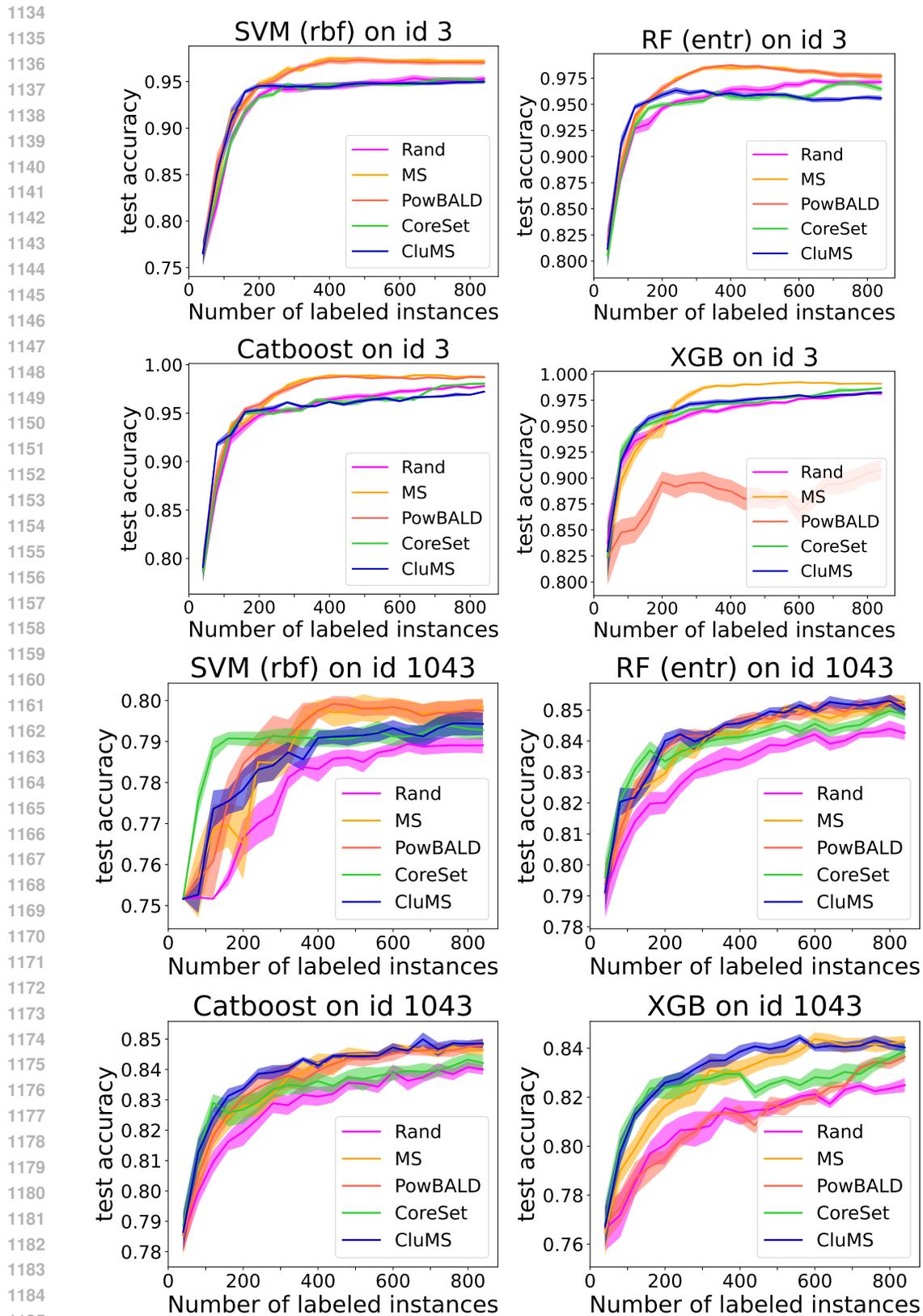


Figure 8: Budget curves for different ALPs on the dataset with OpenML ID 3 and 1043, considering the **large** setting.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

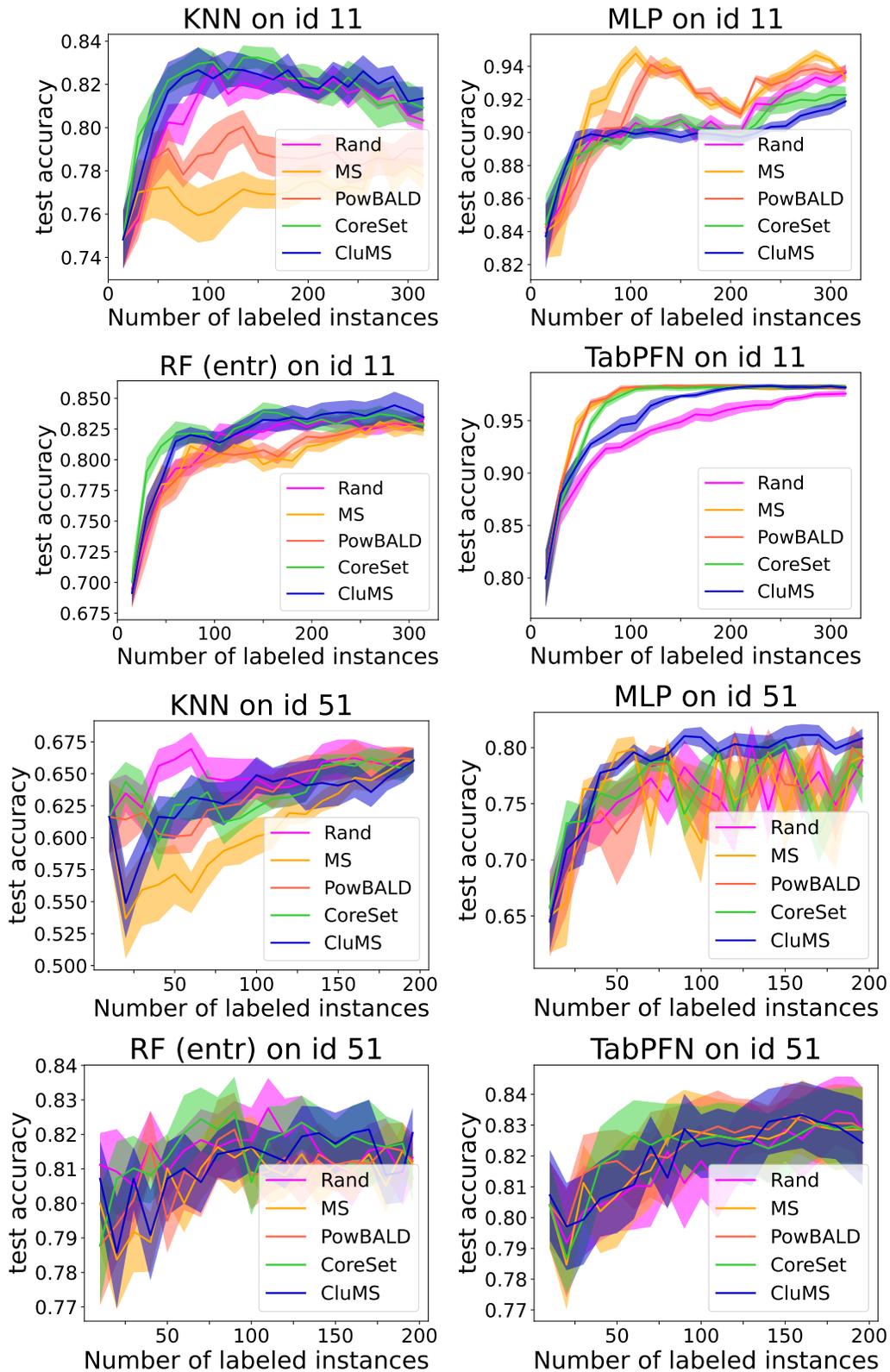


Figure 9: Budget curves for different ALPs on the dataset with OpenML ID 11 and 51, considering the **small** setting.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

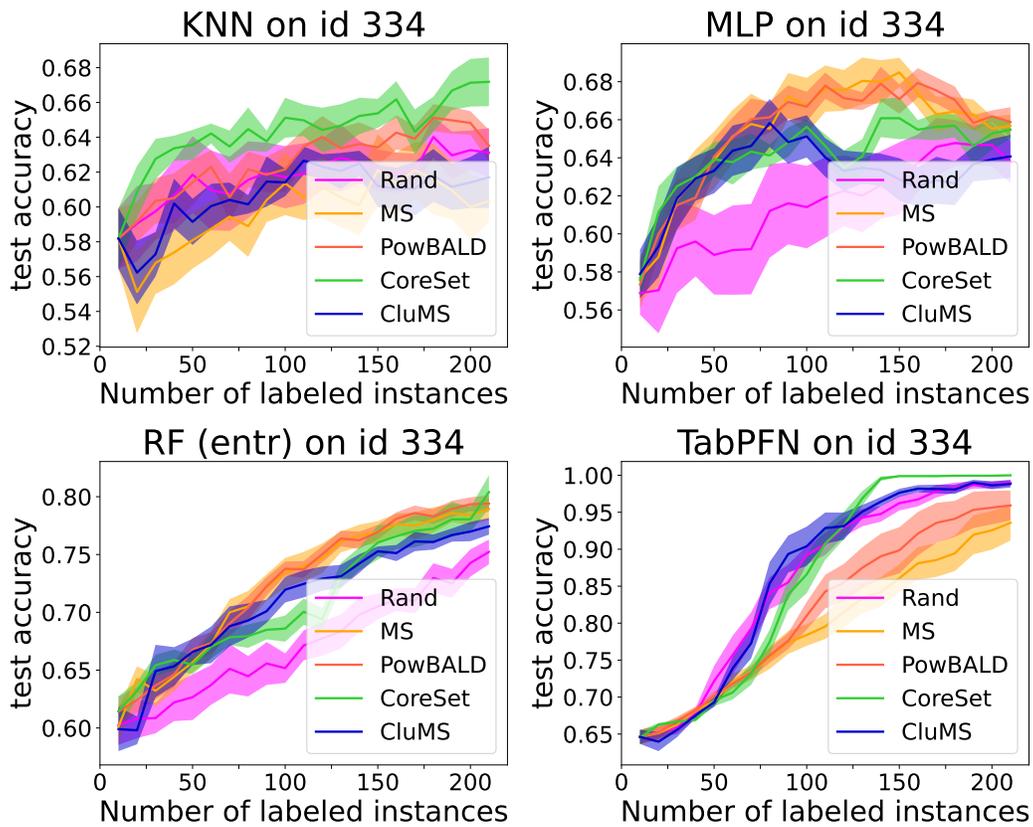


Figure 10: Budget curves for different ALPs on the dataset with OpenML ID 334, considering the **small** setting.