

1 Widespread use of invalid statistical tests in biomedical machine learning

2
3 Tianchu Zeng^{1,2,3,4,5,6,7*}, Hetu Li^{1,3,4,5,6,7*}, Shaoshi Zhang^{1,2,3,4,5,6,7,8*}, Yan Quan Tan^{1,2,3,4,5,6,7},
4 Fang Tian^{1,2,3,4,5,6,7}, Csaba Orban^{1,3,4,5,6,7}, Lijun An⁹, Wanyu Che^{1,3,4,5,6,7}, Jingwen
5 Cheng^{1,3,4,5,6,7,8}, Joanna Su Xian Chong^{1,4,5,6}, Niousha Dehestani^{1,4,5,6}, Zijian Dong^{1,2,4,5,6}, Xin
6 Li^{1,3,4,5,6,7}, Zhizhou Li^{1,4,5,6}, Mervyn Jun Rui Lim^{1,4,5,6,7,10}, Yi Lin^{1,4,5,6}, Qinrui Ling¹¹, Zijie
7 Ling^{1,3,4,5,6,7}, Xi Zhi Low^{1,4,5,6}, Sina Mansour L.^{1,4,5,6,7,12}, Eric Kwun Kei Ng^{1,4,5,6}, Thuan Tinh
8 Nguyen^{1,4,5,6}, Leon Qi Rong Ooi^{1,2,3,4,5,6,7,8}, Shreya Pande^{1,2,3,4,5,6,7,8}, Xing Qian^{1,4,5,6}, Jingxuan
9 Ruan^{1,4,5,6}, Ziwen Wang^{1,3,4,5,6,7}, Yapei Xie^{1,3,4,5,6,7}, Chen Zhang^{1,2,3,4,5,6,7}, Yichi Zhang^{1,4,5,6},
10 Kaustubh Patil^{13,14}, Linden Parkes¹⁵, Elvisha Dhamala^{16,17}, Sidhant Chopra^{18,19}, Andrew
11 Zalesky^{12,20}, Avram Holmes^{21,22}, Simon Eickhoff^{13,14}, Juan Helen Zhou^{1,2,4,5,6,7,8}, Olivier
12 Renaud²³, Nico Dosenbach^{24,25,26,27,28,29}, Konrad Kording^{30,31}, Danilo Bzdok^{32,33,34}, Thomas E.
13 Nichols^{35,36+}, B.T. Thomas Yeo^{1,2,3,4,5,6,7,8,37+}

14
15 ¹Centre for Sleep & Cognition & Centre for Translational Magnetic Resonance Research,
16 Yong Loo Lin School of Medicine, National University of Singapore, Singapore;
17 ²Department of Electrical and Computer Engineering, National University of Singapore,
18 Singapore; ³N.1 Institute for Health, National University of Singapore, Singapore;
19 ⁴Department of Medicine, Yong Loo Lin School of Medicine, National University of
20 Singapore, Singapore; ⁵Healthy Longevity Translational Research Programme, Yong Loo Lin
21 School of Medicine, National University of Singapore, Singapore; ⁶Human Potential
22 Translational Research Programme, Yong Loo Lin School of Medicine, National University
23 of Singapore, Singapore; ⁷Institute for Digital Medicine (WisDM), Yong Loo Lin School of
24 Medicine, National University of Singapore, Singapore; ⁸Integrative Sciences and
25 Engineering Programme (ISEP), National University of Singapore, Singapore; ⁹Department
26 of Clinical Sciences Malmö, SciLifeLab, Lund University, Lund, Sweden; ¹⁰Division of
27 Neurosurgery, Department of Surgery, National University Hospital, Singapore, Singapore;
28 ¹¹Department of Electronic Engineering and Information Science, University of Science and
29 Technology of China, Hefei, China; ¹²Systems Neuroscience Lab, Department of Psychiatry,
30 The University of Melbourne, Parkville, Victoria, Australia; ¹³Institute of Neuroscience and
31 Medicine, Brain & Behaviour (INM-7), Research Center Jülich, Jülich, Germany; ¹⁴Institute
32 for Systems Neuroscience, Medical Faculty, Heinrich-Heine University Düsseldorf,
33 Düsseldorf, Germany; ¹⁵Department of Psychiatry, Brain Health Institute, Rutgers University,
34 Piscataway, NJ, USA; ¹⁶Institute of Behavioral Sciences, Feinstein Institutes for Medical
35 Research, Manhasset, USA; ¹⁷Department of Psychiatry, Donald and Barbara Zucker School
36 of Medicine at Hofstra/Northwell, Uniondale, Hempstead, USA; ¹⁸Orygen, The National
37 Centre of Excellence in Youth Mental Health, Melbourne, Victoria, Australia; ¹⁹Center for
38 Youth Mental Health, University of Melbourne, Melbourne, Victoria, Australia;
39 ²⁰Department of Biomedical Engineering, The University of Melbourne, Parkville, Victoria,
40 Australia; ²¹Department of Psychiatry, Rutgers University, Piscataway, NJ, USA; ²²Center
41 for Advanced Human Brain Imaging Research, Rutgers University, Piscataway, NJ, USA;
42 ²³Department of Psychology, University of Geneva, Genève, Switzerland; ²⁴Mallinckrodt
43 Institute of Radiology, Washington University School of Medicine, St Louis, MO, USA;
44 ²⁵Allied Labs for Imaging Guided Neurotherapies (ALIGN), Washington University School
45 of Medicine, St Louis, MO, USA; ²⁶Department of Neurology, Washington University

46 School of Medicine, St Louis, MO, USA; ²⁷Department of Paediatrics, Washington
47 University School of Medicine, St Louis, MO, USA; ²⁸Department of Biomedical
48 Engineering, Washington University, St Louis, MO, USA; ²⁹Department of Psychological
49 and Brain Sciences, Washington University, St Louis, MO, USA; ³⁰CIFAR Learning in
50 Machines and Brains program, Toronto, Ontario, Canada; ³¹Departments of Bioengineering
51 and Neuroscience, University of Pennsylvania, Philadelphia, PA, USA; ³²The Neuro,
52 McConnell Brain Imaging Centre, Department of Biomedical Engineering, Montreal,
53 Quebec, Canada; ³³Faculty of Medicine, School of Computer Science, McGill University,
54 Montreal, Quebec, Canada; ³⁴Mila–Quebec Artificial Intelligence Institute, Montreal,
55 Quebec, Canada; ³⁵Big Data Institute, Li Ka Shing Centre for Health Information and
56 Discovery, Nuffield Department of Population Health, University of Oxford, Oxford, UK;
57 ³⁶Centre for Integrative Neuroimaging (OxCIN), FMRIB, Nuffield Department of Clinical
58 Neurosciences, University of Oxford, Oxford, UK; ³⁷Athinoula A. Martinos Center for
59 Biomedical Imaging, Massachusetts General Hospital, Charlestown, USA

60

61

62 * equal contributions, + equal contributions

63

64 **Address correspondence to:**

65

66 B.T. Thomas Yeo
67 Yong Loo Lin School of Medicine
68 National University of Singapore
69 Email: thomas.yeo@nus.edu.sg

70

71 Thomas E Nichols
72 Big Data Institute, OxCIN
73 University of Oxford
74 Email: thomas.nichols@bdi.ox.ac.uk

75

76

77

Abstract

78 Machine learning is accelerating biomedical research. Cross-validation is widely used to
79 compare predictive performance – not only to benchmark algorithms, but also to inform
80 scientific applications, such as ranking biomarkers. However, prediction performance
81 estimates across cross-validation folds are not independent. Standard tests for comparing
82 prediction performance (e.g., paired t-test) assume independence and can therefore inflate
83 false positive rates. In a PRISMA-guided meta-analysis of 210 studies (impact factor ≥ 15 , 1
84 June 2020 – 1 June 2025), we find that 97% ignored fold dependence when comparing
85 prediction performance. This problem is ubiquitous across scientific fields and unaffected by
86 impact factor, rigor-promoting policies, or open science practices. Simulations across 420
87 scenarios spanning four diverse datasets show that ignoring fold dependence leads to invalid
88 false positive control in most settings. Repeated cross-validation further compounds this
89 problem, with false positive rates rising toward 100% as the number of repetitions grows.
90 Existing fold-dependence-aware tests rely on strong assumptions because the variance of
91 fold-level statistics and the between-fold correlation cannot be disentangled under standard
92 cross-validation. We therefore propose the SHARP (Split-HAlf RePeated) test, a simple
93 modification to standard cross-validation that enables direct estimation of variance and
94 correlation. Benchmarked against 12 tests, SHARP provides the best overall balance of false-
95 positive control, statistical power, and confidence-interval calibration across simulation
96 schemes. We conclude by providing best practices and reporting guidelines for valid model
97 comparison inference in biomedical machine learning and beyond.

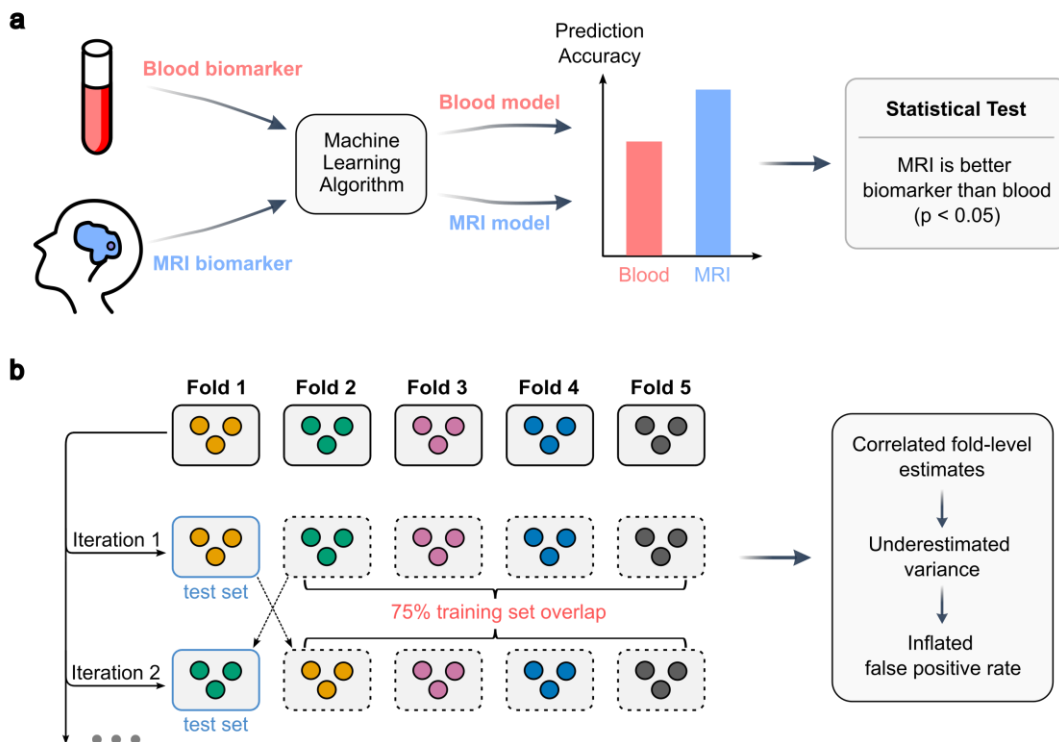
98

1 Introduction

99 Biomedical research increasingly relies on machine learning to extract predictive information
100 from complex, high-dimensional data (Rajpurkar et al., 2022; Moor et al., 2023; Bzdok et al.,
101 2024; Perez-Lopez et al., 2024). As algorithms and data modalities proliferate, empirical
102 comparisons of predictive performance have become routine, guiding decisions about which
103 algorithms to deploy or which biomarkers to prioritize (Greene et al., 2018; Wagner et al.,
104 2023; Carrasco-Zanini et al., 2024; Yoo et al., 2025). Consider a hypothetical study in which
105 the same algorithm is trained to predict dementia progression from either blood or MRI
106 biomarkers (Fig. 1a). Because the algorithm is held constant, superior performance of the
107 MRI-based model would indicate that MRI carries more information about dementia
108 progression than blood. Here, the goal extends beyond benchmarking algorithms to informing
109 scientific and clinical applications.

110

111 A standard approach for comparing predictive performance is K-fold cross-validation. The
112 dataset is partitioned into K non-overlapping subsets (“folds”) of observations, such as brain
113 scans or blood samples. In each iteration, one fold serves as the test set while the remaining
114 $K - 1$ folds form the training set (Fig. 1b). When applied to two competing biomarkers (or
115 algorithms), this procedure yields K pairs of prediction performance estimates that can be
116 compared using a statistical test. Although test folds do not overlap across iterations, training
117 sets overlap substantially, and each test fold contributes to the training set in all other
118 iterations (Fig. 1b). The prediction performance estimates are therefore statistically
119 dependent across folds (Dietterich, 1998; Bates et al., 2024). Conventional statistical tests
120 (e.g., paired t-tests or Wilcoxon signed-rank tests) assume independence, so they
121 underestimate variability and inflate false positive rates (Fig. 1b; Nadeau & Bengio, 2003;
122 Jafrasteh et al., 2025). Repeated cross-validation, often recommended to improve stability
123 (Bouckaert & Frank, 2004; Varoquaux et al., 2017), introduces further overlap across training
124 and test sets, which can worsen the problem (Nadeau & Bengio, 2003).



125

126

127

Figure 1. Comparing the predictive power of biomarkers using machine learning and statistical dependence induced by cross-validation. a. Machine learning model comparison

128 can inform scientific and clinical applications. In this hypothetical example, the same
129 algorithm is trained to predict dementia progression from blood or MRI biomarkers, and a
130 statistical test compares their performance on held-out data. Because the algorithm is held
131 constant, superior performance of the MRI model would suggest that MRI is more
132 informative than blood. **b.** *K*-fold cross-validation (CV), illustrated for $K = 5$. The dataset is
133 partitioned into five folds. In each iteration, one fold serves as the test set (blue outline) and
134 the remaining four folds form the training set (dashed outline). Training sets overlap
135 substantially across iterations, and each test fold contributes to the training set in all other
136 iterations. This induces positive correlations across fold-level estimates. Treating the fold-
137 level estimates as independent underestimates the true variance, which inflates false positive
138 rates (Nadeau & Bengio, 2003). Methods Section 4.1 discusses the properties of other
139 variants of standard cross-validation.

140
141

142 Although fold-dependence has been recognized for decades (Dietterich, 1998; Nadeau &
143 Bengio, 2003), its prevalence and consequences remain unclear. We ourselves have
144 inadvertently applied statistical tests that ignore fold dependence (Yeo et al., 2010; Mansour
145 L. et al., 2021; Chopra et al., 2024), raising a broader question: how often is fold dependence
146 ignored in biomedical research, and does this vary across scientific fields? Furthermore, if
147 existing publishing safeguards – editorial selectivity, statistical review, methodological
148 reporting standards, and the scrutiny enabled by open data and code – were catching invalid
149 tests, prevalence should be lower in studies meeting these criteria.

150

151 Beyond prevalence, fully addressing the fold-dependence problem requires confronting a
152 fundamental statistical ambiguity. A single run of standard *K*-fold cross-validation yields *K*
153 fold-level estimates of prediction performance differences, and three unknowns: true mean,
154 true variance, and true between-fold correlation. The *K* fold-level differences, however,
155 provide only two sources of information – their sample mean and sample variance. The
156 sample variance underestimates the true variance, and is on average equal to the true variance
157 $\times (1 - \text{true correlation})$. Thus the true variance and correlation cannot be disentangled without
158 additional assumptions (Nadeau & Bengio, 2003; Bengio & Grandvalet, 2004). Existing tests
159 cope by making strong (implicit or explicit) assumptions about the true variance or
160 correlation rather than estimating both quantities from data. When those assumptions fail, the
161 result can be false positive inflation or lower statistical power.

162

163 Here we perform a PRISMA-guided meta-analysis of 210 PubMed-listed studies (impact
164 factor ≥ 15), quantifying how often fold dependence is ignored, and evaluating whether
165 prevalence varies by scientific field, impact factor, journal policy, or open science practice.
166 We then develop and apply a battery of 420 simulation scenarios spanning image recognition,
167 neuroimaging, ecology, and systems biology. Invalid statistical tests fail to control false
168 positives in most settings, with false positive rates rising toward 100% under repeated cross-
169 validation. We propose the SHARP (Split-Half Repeated) test, which modifies standard
170 cross-validation to generate pairs of independent statistics, enabling direct estimation of
171 variance and correlation. Benchmarking against 12 existing tests, SHARP reliably controls
172 false positives and yields well-calibrated confidence intervals, while matching or exceeding
173 the statistical power of valid tests. Finally, our meta-analysis also reveals that over half the
174 studies comparing model performance either did not apply a statistical test or report
175 confidence intervals, or described their procedure too vaguely to evaluate, so we also provide
176 a concrete set of reporting recommendations for predictive model comparison.

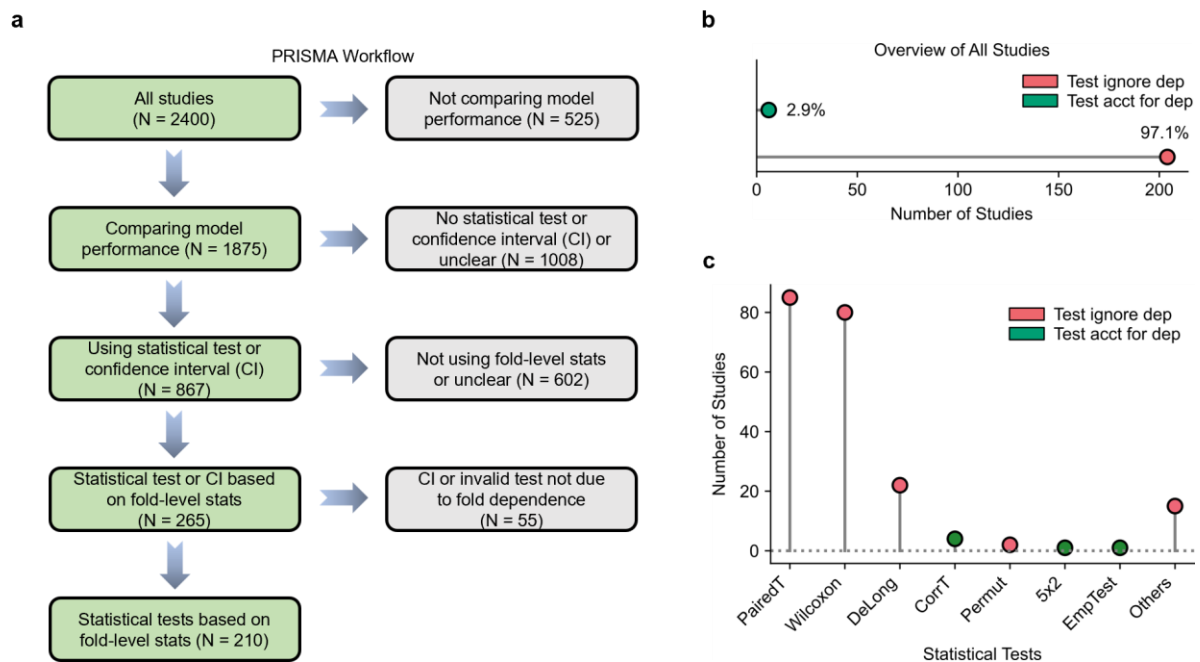
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196

2 Results

2.1 Invalid statistical tests in 97% of studies

To assess how often invalid statistical tests are used, we conducted a meta-analysis following the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines (Page et al., 2021; Fig. 2a; Methods Section 4.2). Using PubMed, we identified original research articles published between 1 June 2020 and 1 June 2025 that used cross-validation to compare predictive performance. We restricted inclusion to journals with impact factor ≥ 15 , which typically have the most stringent reporting standards, statistical review, and transparency policies. This is a conservative criterion: if the problem is prevalent even in venues with the strongest methodological safeguards, it is likely to be as prevalent elsewhere.

The initial search yielded 2,400 studies, of which 1,875 compared model performance (Fig. 2a). Of these, 1,008 studies (54%) did not clearly apply a statistical test or report confidence intervals when comparing model performance. In some cases, models were compared based on point estimates alone; in others, reporting was too vague to determine whether any statistical test was performed. We focused on the remaining 867 studies that clearly applied a statistical test or reported confidence intervals (Fig. 2a). Of these, 602 were excluded because the actual statistical procedure was unclear or did not use fold-level statistics (Methods Section 4.2.2), resulting in 265 studies (Fig. 2a).



197
198
199
200
201
202
203
204
205
206
207
208
209

Figure 2. Widespread use of invalid statistical tests in biomedical studies using cross-validation to compare predictive performance. **a.** PRISMA flow diagram of the meta-analysis. A PubMed search (1 June 2020–1 June 2025; journals with impact factor ≥ 15) identified 2,400 records, of which 1,875 compared model performance. Among these, 867 clearly applied statistical tests or confidence intervals. We further excluded 602 studies without fold-level statistics or with unclear procedures, one study with a permutation test invalid for reasons unrelated to fold dependence, and 54 studies reporting only confidence intervals, leaving 210 studies for analysis. **b.** Breakdown of the 210 studies by statistical approach: tests assuming fold independence (red) versus tests accounting for fold dependence (green). **c.** Breakdown of specific statistical tests across the 210 studies. Red indicates tests that assume fold independence, i.e., invalid tests. Green indicates tests that account for fold

210 dependence. “PairedT” refers to the uncorrected resampled paired t-test (Nadeau & Bengio,
211 2003); “Wilcoxon” refers to either the Wilcoxon signed-rank test (Wilcoxon, 1945) or the
212 Wilcoxon rank-sum test (Mann & Whitney, 1947); “DeLong” refers to DeLong’s test
213 (DeLong et al., 1988); “Permut” refers to paired permutation test (Edgington & Onghena,
214 2007); “CorrT” refers to the corrected resampled paired t-test (Nadeau & Bengio, 2003);
215 “5×2” refers to the “5×2 paired t-test” (Dietterich, 1998); and “EmpTest” refers to a variant
216 of the empirical test of differences (Parkes et al., 2021a).

217
218

219 We excluded one study using a permutation test that was invalid for reasons unrelated to fold
220 dependence (see Supplementary Results). We also excluded studies reporting only
221 confidence intervals, because the method of computation was rarely reported with sufficient
222 clarity to assess validity (see Methods Section 4.2). This resulted in a final set of 210 studies
223 (Fig. 2a).

224

225 The winnowing from 1,875 to 210 reflects broader reporting deficiencies in the ML literature
226 — from studies where it was unclear whether any statistical test was performed, to studies
227 with insufficiently described statistical procedures. These deficiencies are outside the scope
228 of our fold-dependence analysis but warrant attention in their own right (see Discussion
229 Section 3.4). Our meta-analysis focuses specifically on the 210 studies with sufficient clarity
230 to assess fold dependence.

231

232 Of the 210 studies (Fig. 2b), 204 (97%) used statistical tests that assumed independent fold-
233 level estimates. The paired t-test and Wilcoxon signed-rank test were the most common
234 invalid tests (Fig. 2c). Further breakdown of meta-analysis results is found in Supplementary
235 Results and supporting quotations from each study are provided in Supplementary File 1. As
236 we show below (Section 2.4), these tests inflate FPR to 19% on average – nearly 4× the
237 nominal 5% level – under a single round of 10-fold cross-validation, and break down entirely
238 (FPR → 100%) under repeated cross-validation, a commonly recommended practice.

239

240

241 **2.2 Widespread invalid statistical tests over time and scientific fields**

242 We next examined whether the prevalence of invalid fold-based statistical tests (among the
243 210 studies) showed any discernible pattern over time or across scientific fields. Here,
244 scientific fields are operationalized as Web of Science Journal Citation Reports subject
245 categories (Methods Section 4.2.3).

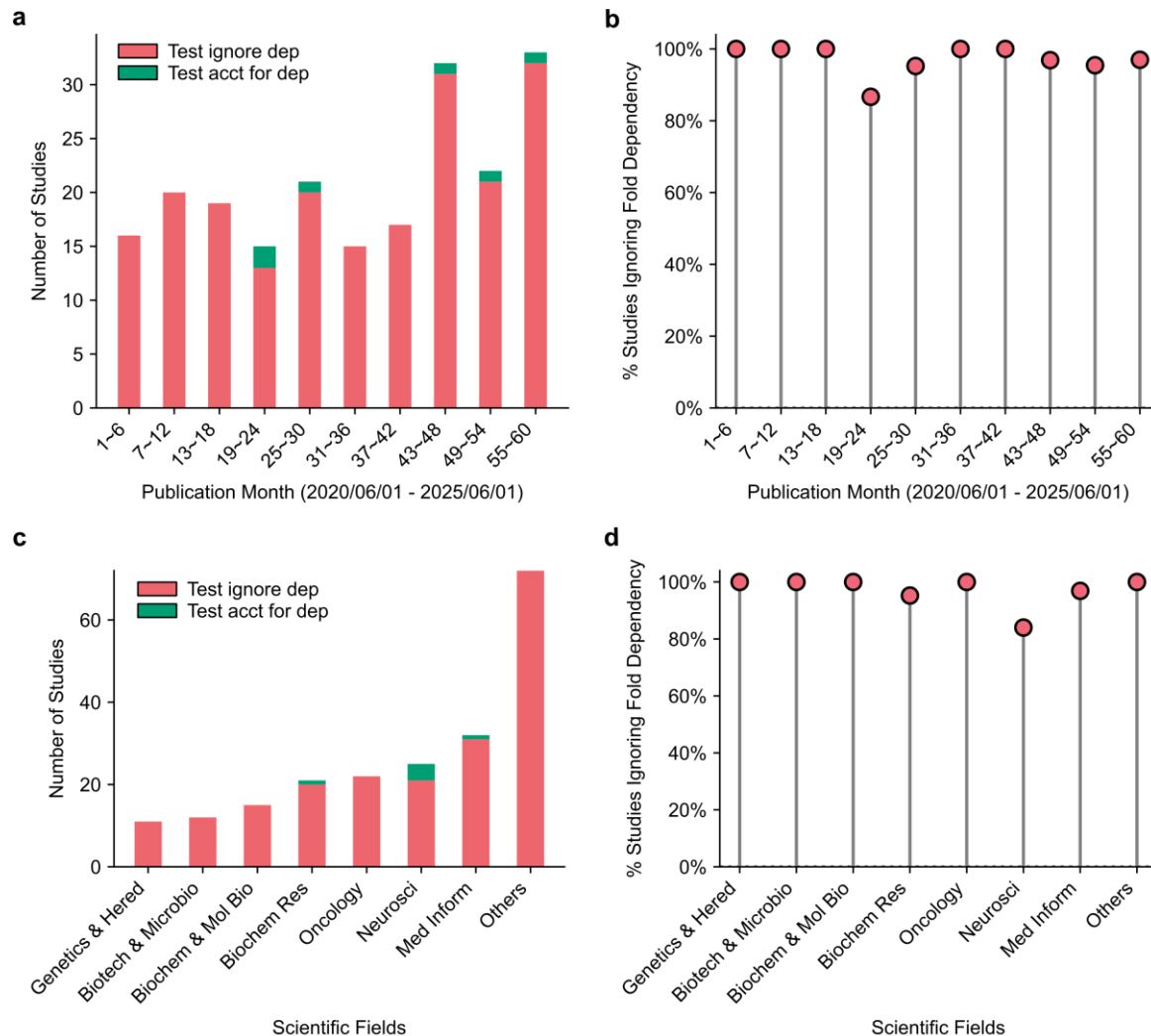
246

247 The number of studies grew by 14.4% per year from 1 June 2020 to 1 June 2025 (Poisson
248 regression $p = 6e-3$; Fig. 3a), reflecting the increasing use of cross-validation to compare
249 machine learning models. However, the proportion of studies ignoring fold dependence
250 showed no detectable change over time, remaining persistently high across all periods
251 (permutation test $p = 0.63$; Fig. 3b).

252

253 Across scientific fields, reliance on invalid fold-independence assumptions was ubiquitous
254 (Figs. 3c,d). Among fields with at least 10 studies, 100% of studies in “Genetics & Heredity”,
255 “Biotechnology & Applied Microbiology”, and “Biochemistry & Molecular Biology”
256 ignored fold dependence (Fig. 3d). Neuroscience had the lowest proportion, at 84.0%.
257 Nevertheless, the proportion did not differ significantly across fields (permutation test $p =$
258 0.07). Together, these results show that the use of invalid tests has neither declined over time
259 nor varied meaningfully across scientific fields.

260



261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

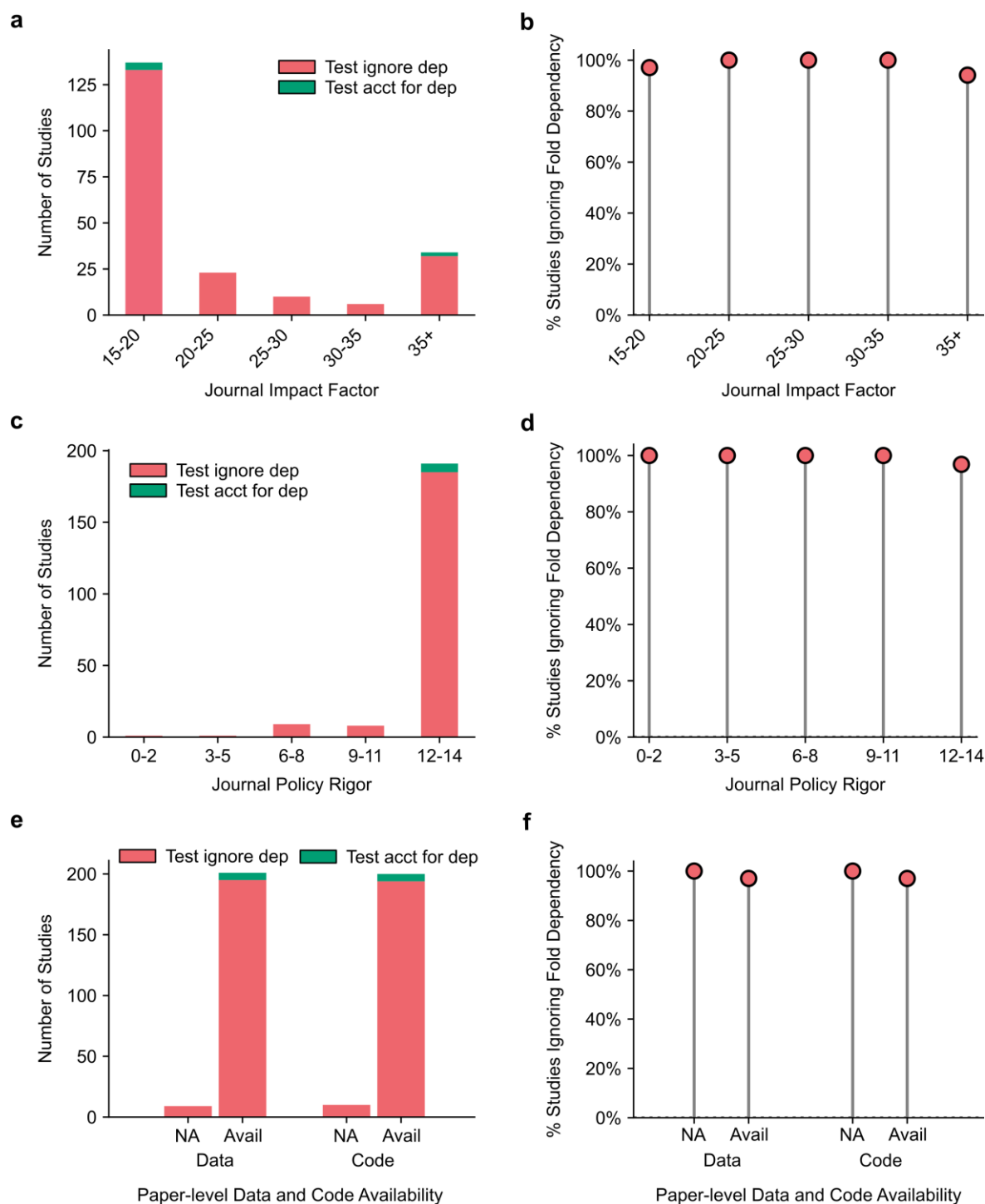
282

Figure 3. Persistent risk of inflated false positives from fold dependence over time and across scientific fields. **a.** Number of studies by publication period, in 6-month intervals from 1 June 2020 to 1 June 2025 (N = 210). **b.** Proportion of studies ignoring fold dependence over time. The proportion showed no detectable change (permutation test, $p = 0.63$). **c.** Number of studies by scientific field. The 210 studies spanned 33 fields. Fields with at least 10 studies are shown separately, while the remaining 26 fields are grouped as “Others”. **d.** Proportion of studies ignoring fold dependence by scientific field. The proportion of studies ignoring fold dependence did not differ significantly across scientific fields (permutation test, $p = 0.07$). Note that the “Others” category was excluded from the statistical comparison across fields. Details about statistical tests are found in Methods Section 4.2.7.

2.3 Invalid tests used regardless of impact factor, policy or open science practices

We next examined whether the use of invalid statistical tests was associated with journal impact factor, journal policies promoting scientific rigor, or open science practices. If existing publishing safeguards were filtering out invalid tests, we would expect lower prevalence in journals or studies satisfying these criteria — through editorial selectivity, explicit methodological policies, or the post-hoc scrutiny enabled by shared data and code.

283 First, we stratified studies by the impact factor of the publishing journal (Fig. 4a). The
284 proportion of studies ignoring fold dependence did not vary across impact factor bins
285 (permutation test $p = 0.59$; Fig. 4b).
286



287
288

289 **Figure 4. Neglect of fold-dependence persists across impact factor, journal policies for**
290 **scientific rigor and open science practices. a.** Distribution of studies by journal impact
291 factor. **b.** Proportion of studies ignoring fold dependence by impact factor bins, calculated as
292 the number of studies using statistical tests ignoring fold dependence divided by the total
293 number of studies in each bin. No trend was detected (permutation test $p = 0.59$). **c.**
294 Distribution of studies by journal policy rigor score, ranging from 0 (least rigorous) to 14

295 (most rigorous). **d.** Proportion of studies ignoring fold dependence by journal policy rigor
296 score. Increasing rigor did not reduce the use of invalid statistical tests (permutation test $p =$
297 0.86). **e.** Number of studies ignoring fold dependence stratified by data and code availability.
298 Most studies provided data, code or both. **f.** Proportion of studies ignoring fold dependence
299 stratified by data and code availability. The proportion was not significantly associated with
300 data or code availability (permutation test $p = 1.00$ for both).

301
302
303

304 For each journal, we assigned a rigor score (0 to 14) based on policies for transparency and
305 reproducibility (e.g., code and data availability) and policies for statistical and
306 methodological guidance (e.g., reporting checklists). Most journals had relatively stringent
307 policies (Fig. 4c). Nevertheless, the proportion of studies ignoring fold dependence remained
308 uniformly high regardless of journal rigor score (permutation test $p = 0.86$; Fig. 4d).

309

310 Finally, most studies provided data, code, or both (Fig. 4e). Among studies that did not
311 provide data, code, or both ($N = 16$), 100% ignored fold dependence (Fig. 4f). However, the
312 proportion of studies ignoring fold dependence was not associated with data availability
313 (permutation test $p = 1.00$) or code availability (permutation test $p = 1.00$).

314

315 Together with Section 2.2, these results show that the prevalence of invalid tests is not
316 explained by year, field, journal impact factor, journal rigor policies, or open science
317 practices, pointing instead to a systemic gap in fold dependence awareness (see Discussion
318 Section 3.3).

319

320

321 **2.4 Ignoring fold dependence leads to poor control of false positive rate (FPR)**

322 To assess the impact of invalid tests on false positive rates (FPR), we analyzed four datasets
323 spanning image recognition (EMNIST Digits; Cohen et al., 2017), neuroimaging (UK
324 Biobank; Alfaro-Almagro et al., 2018), ecological classification (Coverttype; Blackard, 1998),
325 and systems biology (KEGG Metabolic Pathway; Muhammad Naeem, 2011). Further details
326 about the datasets are provided in Methods Section 4.3.

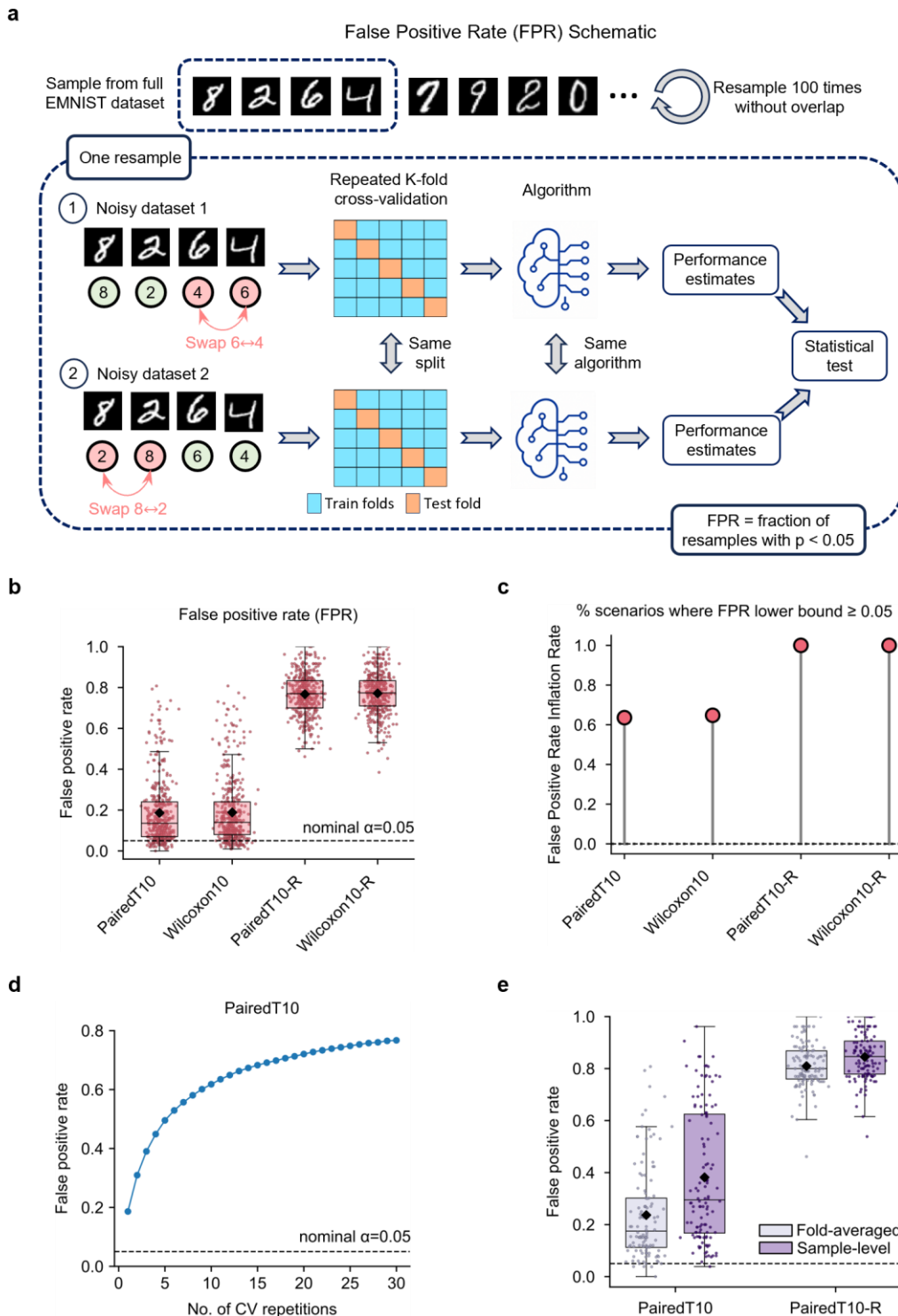
327

328 To estimate FPR, we constructed pairs of trained models with identical expected predictive
329 performance and compared them using a statistical test (Fig. 5a). First, we randomly sampled
330 a dataset of a given sample size (e.g., $N = 1,000$) from one of the four datasets. Next, we
331 generated two noisy versions of this dataset by independently permuting the target labels
332 (e.g., digit labels for EMNIST) for the same fixed fraction of samples (e.g., 10%) in each
333 copy. Finally, the same machine learning algorithm and cross-validation scheme were applied
334 to both noisy copies. The independent noise ensured that the two trained models were
335 distinct, while the matched noise level ensured equal expected performance under the null.
336 Therefore, a well-calibrated test at significance threshold $\alpha = 0.05$ should reject the null
337 hypothesis no more than 5% of the time (Fig. 5a).

338

339 For each combination of dataset, sample size, and noise level — hereafter a "condition" —
340 we repeated the above procedure on non-overlapping subsamples of the full dataset (typically
341 100 repetitions per condition; see Methods Section 4.4), yielding one p-value per repetition.
342 The FPR for each condition was estimated as the fraction of repetitions in which the null was
343 rejected. The independent repetition also enabled 95% confidence intervals for the FPR. For
344 each dataset, we varied the noise level (10%, 20%, ..., 100% labels permuted) and sample

345 size ($N = 100, 500, 1,000, 2,000$; two datasets each excluded one sample size; see Methods
 346 Section 4.4), yielding $(4+4+3+3) \times 10 = 140$ conditions. For each condition, we considered
 347 three hyperparameter-selection strategies: fixed hyperparameters, a single training-validation
 348 split, and nested cross-validation. This yielded $140 \times 3 = 420$ scenarios.



349 **Figure 5. Ignoring fold dependence leads to poor control of false positive rates.** **a.** False
 350 positive rate (FPR) simulation illustrated for the EMNIST dataset. We randomly sampled a
 351 dataset (e.g., $N = 1,000$) from a full dataset, then generated two noisy versions by
 352 independently permuting the target labels (ground-truth digits) for a fixed percentage of
 353 samples (e.g., 10%). For example, noisy dataset 1 (top) has the labels for digits 6 and 4
 354

355 swapped, while noisy dataset 2 (bottom) has the labels for digits 8 and 2 swapped, reflecting
356 the independence of the two permutations. Using identical cross-validation splits and the
357 same machine learning algorithm on each noisy dataset, we obtained pairs of cross-validated
358 performance estimates — e.g., 300 pairs for 10-fold cross-validation repeated 30 times. A
359 statistical test was applied to the paired vector to obtain a p-value. Because the two noisy
360 datasets had the same noise level, the two trained models had identical expected predictive
361 performance. At $\alpha = 0.05$, a well-calibrated test should reject no more than 5% of the time.
362 This procedure was repeated up to 100 times using non-overlapping subsamples. **b.** FPR
363 across four datasets and four sample sizes. Each boxplot shows 420 datapoints, one per
364 scenario. The black dot indicates the mean FPR. The dashed line indicates the nominal FPR
365 (0.05). Abbreviations: PairedT10 / Wilcoxon10: paired t-test / Wilcoxon signed-rank test
366 based on 10-fold cross-validation; PairedT10-R / Wilcoxon10-R: same tests with 10-fold
367 cross-validation repeated 30 times. **c.** FPR inflation rate: percentage of the 420 scenarios in
368 which the lower bound of the 95% confidence interval for FPR exceeded 0.05. **d.** FPR of the
369 paired t-test as a function of the number of 10-fold cross-validation repetitions, averaged over
370 420 scenarios. The dashed line indicates the nominal FPR (0.05). Wilcoxon signed-rank test
371 yields the same conclusions (not shown). **e.** FPR of the paired t-test using sample-level versus
372 fold-averaged statistics, across 120 scenarios in the KEGG Metabolic Pathway dataset (one
373 datapoint per scenario). The black dot indicates the mean FPR. The dashed line indicates the
374 nominal FPR (0.05).

375
376

377 Across the 420 scenarios, the two most widely used tests (the paired t-test and Wilcoxon
378 signed-rank test; Fig. 2c) exhibited an elevated FPR of ~19% when 10-fold cross-validation
379 was performed once (Fig. 5b). Type I error was not controlled (lower bound of the 95% CI
380 exceeded 0.05) in ~64% of scenarios, hereafter the “FPR inflation rate” (Fig. 5c). Moreover,
381 when 10-fold cross-validation was repeated 30 times – as is often recommended for stability
382 (Bouckaert & Frank, 2004; Varoquaux et al., 2017) – the FPR jumped to 77% (Fig. 5b) and
383 the tests were invalid in 100% of scenarios. FPR rose monotonically with the number of
384 repetitions (Fig. 5d), eventually reaching 100%. We emphasize that the problem lies not with
385 repeated cross-validation itself, which has well-established benefits, but with its interaction
386 with tests that assume fold independence.

387

388 The escalation of FPR with repeated cross-validation is not specific to our simulations: any
389 test that assumes fold independence will behave this way when repeated cross-validation is
390 applied. Even when two models have the same expected performance in the population (i.e.,
391 null hypothesis is true), a finite dataset will almost certainly show a small, non-zero sample-
392 specific difference between the two models. Standard tests such as the paired t-test treat each
393 fold-level observation as independent, effectively assuming the sample size grows with the
394 number of folds and repetitions. Since any effect, however small, becomes significant when
395 the sample size is large enough, the FPR is expected to approach 100% as the number of
396 repetitions grows.

397

398

399 **2.5 Sample-level statistics compound false positive rate inflation**

400 Section 2.4 focused on tests using “fold-averaged” statistics. Our meta-analysis revealed that
401 some studies instead used “sample-level” statistics, which may inflate FPR further. Consider
402 a study comparing two models via 10-fold cross-validation on 1,000 samples. A fold-
403 averaged approach averages performance within each fold, yielding 10 paired metrics for the
404 statistical test. A sample-level approach forgoes this aggregation and enters all 1,000 paired

405 metrics directly into the test. Both approaches derive values from test folds and were
406 therefore included in our meta-analysis. Because statistics within the same fold might be even
407 more strongly correlated than statistics across folds, sample-level statistics may inflate FPR
408 further.

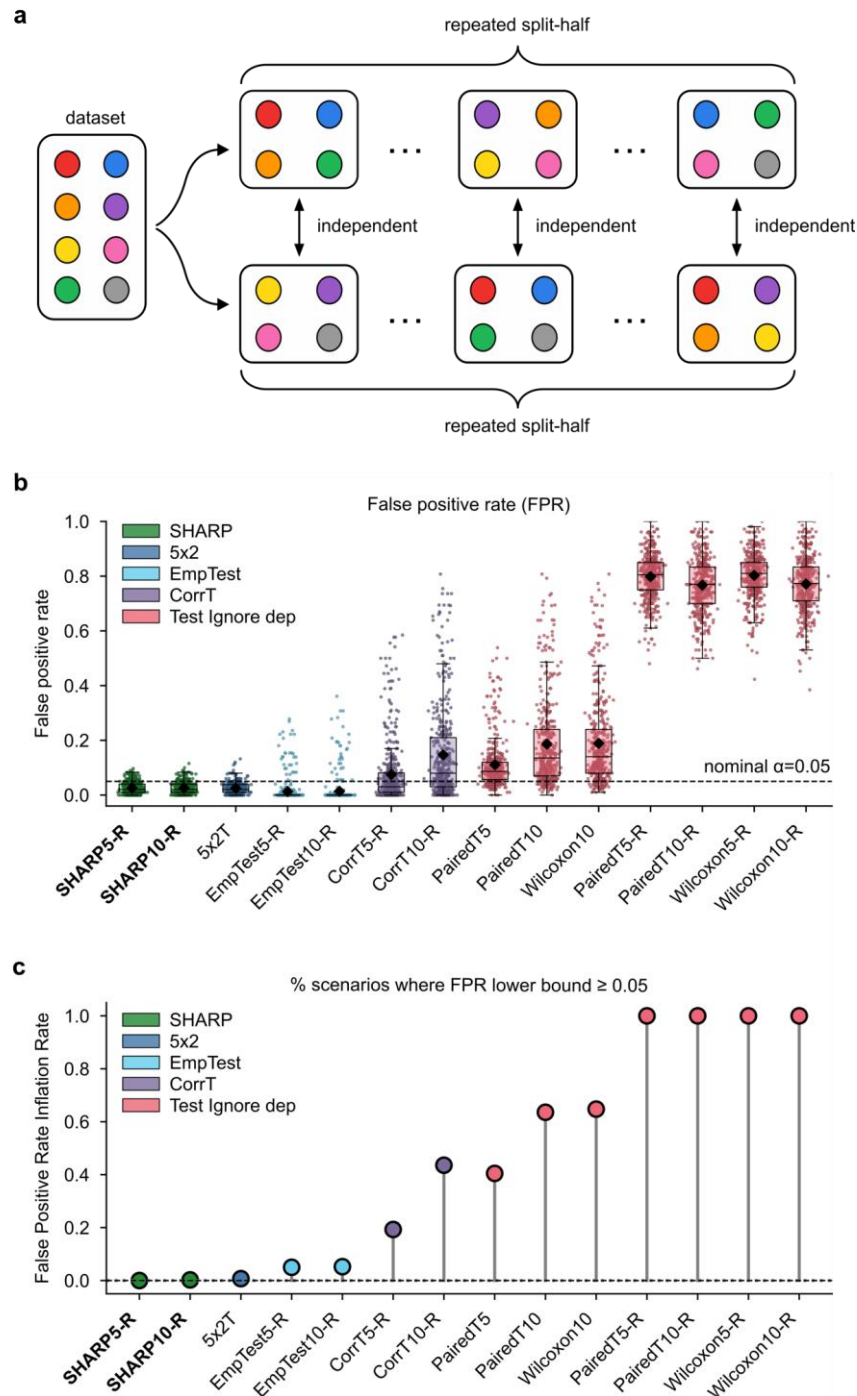
409
410 We first examined this issue using the KEGG Metabolic Pathway dataset. As expected,
411 switching from fold-averaged to sample-level statistics (while still ignoring fold dependence)
412 elevated FPR — from 24% to 38% for a single run of 10-fold cross-validation, and from 81%
413 to 84% when repeated 30 times (Fig. 5e). We next examined DeLong's test, another popular
414 sample-level test (Fig. 2c). Because DeLong's test applies only to binary classification
415 (Methods Section 4.5.4), we evaluated it on the Covertyp dataset, where it exhibited FPRs
416 of 21% and 77% for 10-fold cross-validation with one repetition and 30 repetitions,
417 respectively (Fig. S1). For comparison, the paired t-test on fold-averaged statistics in the
418 same dataset yielded FPRs of 13% and 72% (Fig. S2c). Together, these results show that
419 sample-level statistics compound the FPR inflation already caused by ignoring fold
420 dependence.

421
422
423 **2.6 SHARP avoids strong assumptions in existing fold-dependence-aware tests**
424 Having shown that tests ignoring fold dependence severely inflate FPRs (Fig. 5), we next
425 considered tests that account for it (hereafter “fold-aware” tests). Accounting for fold
426 dependence must address a fundamental statistical ambiguity. In a single run of K -fold cross-
427 validation, the K fold-level differences can be used to compute sample mean and sample
428 variance, but the sampling distribution of the average difference involves three unknown
429 parameters: true mean, true variance, and true between-fold correlation. Since expected
430 sample variance = true variance \times (1 – true correlation), the true variance and correlation
431 cannot be jointly estimated without additional assumptions (Nadeau & Bengio, 2003).
432 Repeated K -fold cross-validation inherits this non-identifiability, and no universal
433 distribution-independent unbiased estimator of the standard error exists (Bengio &
434 Grandvalet, 2004).

435
436 Existing tests address this non-identifiability issue in different ways, each with its own
437 limitations. Among the six studies in our meta-analysis that used fold-aware tests (Fig. 2c):
438 four used the corrected resampled t-test (Nadeau & Bengio, 2003), one used the 5×2 paired t-
439 test (Dietterich, 1998), and one used a variant of the empirical test of differences (Parkes et
440 al., 2021a). The corrected resampled t-test assumes that the between-fold correlation is equal
441 to the fraction of samples held out in each test fold (Supplementary Methods S3). If true
442 correlation is higher, FPR is inflated, and if lower, power is diminished. The 5×2 paired t-test
443 uses 5 repetitions of 2-fold cross-validation to approximately correct for the between-fold
444 correlation. However, it estimates the mean performance difference from a single fold,
445 discarding information from the remaining nine folds at severe cost to power (Supplementary
446 Methods S4). The empirical test of differences avoids modeling between-fold correlation by
447 using the empirical distribution of fold-level differences. It is conservative when true
448 correlation ≤ 0.5 and loses power when true correlation is much lower (Supplementary
449 Methods S5). Overall, these assumptions risk inflating FPR or sacrificing power.

450
451 We therefore propose the Split-HAlf RePeated (SHARP) test. In each of J repetitions,
452 SHARP randomly partitions the dataset into two disjoint halves, A and B, and performs K -
453 fold cross-validation independently within each half (Fig. 6a). Fold-level performance
454 differences are averaged within each half, yielding one statistic per half. Because A and B

455 share no observations, the two statistics from a given repetition are independent, while
 456 statistics from different repetitions remain correlated. This within-repetition independence is
 457 sufficient to estimate the variance of each split-half statistic and the across-repetition
 458 correlation, enabling a valid statistical test for assessing differences between two models
 459 (Methods Section 4.6, Supplementary Methods S7). We next evaluate whether this design
 460 delivers in practice, by comparing SHARP's false positive rate (Section 2.7), statistical power
 461 (Section 2.8), and confidence interval coverage (Section 2.9) against existing tests.
 462



463
 464

465 **Figure 6. SHARP test and false positive rates (FPR) of statistical tests across 420**
 466 **simulation scenarios. a.** Schematic of the SHARP test (see Section 2.6 for motivation). The
 467 procedure performs J random splits of the dataset into disjoint halves (A and B) and runs K -

468 fold cross-validation separately on each. Averaging the fold-level differences within a half
469 produces a single statistic, giving two statistics per repetition. The disjoint sampling makes
470 these two statistics independent within a repetition, even though statistics from different
471 repetitions share data and are therefore correlated. This structure provides the two ingredients
472 needed for a valid test: an estimate of each statistic's variance and an estimate of the across-
473 repetition correlation. **b.** Boxplot of FPR for each statistical test, with one value per scenario
474 ($n = 420$). Black dots mark the mean FPR across scenarios. The dashed horizontal line marks
475 the nominal level of 0.05. **c.** FPR inflation rate: the percentage of the 420 scenarios in which
476 the lower bound of the 95% confidence interval for FPR exceeded 0.05, indicating inadequate
477 control of the Type I error rate. **Test naming conventions.** Numeric suffixes "5" and "10"
478 denote 5-fold and 10-fold cross-validation. The "-R" suffix indicates repeated cross-
479 validation: 5-fold repeated 60 times or 10-fold repeated 30 times, yielding $5 \times 60 = 10 \times 30 =$
480 300 fold-level statistics. For example, "CorrT5-R" denotes the corrected resampled t-test
481 evaluated under 60 repetitions of 5-fold cross-validation, while "PairedT10" denotes the
482 naïve paired t-test under a single run of 10-fold cross-validation. "SHARP5-R" denotes the
483 split-half procedure repeated 60 times with 5-fold cross-validation within each half;
484 "SHARP10-R" denotes 30 repetitions with 10-fold cross-validation within each half. The
485 number of repetitions was chosen such that the FPR of tests accounting for fold dependence
486 had stabilized; tests that ignore fold dependence do not stabilize, with FPR increasing toward
487 1 as repetitions grow (Fig. 5d). "CorrT" denotes the corrected resampled t-test (Nadeau &
488 Bengio, 2003). "EmpTest" denotes the empirical test of differences (Parkes et al., 2021a). For
489 details of all tests, see Methods Sections 4.5 and 4.6.

490
491
492

493 **2.7 SHARP: reliable FPR control through direct variance estimation**

494 Across 420 simulation scenarios (Section 2.4), the corrected resampled t-test, 5×2 paired t-
495 test, and empirical test of differences markedly reduced FPR relative to tests ignoring fold
496 dependence (Fig. 6b). The 5×2 paired t-test and empirical test of differences controlled FPR
497 at the nominal 5% level. The corrected resampled t-test was borderline under 5-fold cross-
498 validation and failed under 10-fold, with FPR significantly exceeding 5% in 19% and 44% of
499 scenarios, respectively (Fig. 6c). These conclusions held across datasets (Fig. S2) and sample
500 sizes (Fig. S3). Because simulation results can depend on the underlying modeling
501 assumptions, we repeated the analyses under an alternative scheme (Supplementary Methods
502 S6) – results were unchanged, except that the corrected resampled t-test now controlled FPR
503 (Fig. S4), indicating that its FPR control depends on the simulation scheme.

504

505 Across the same 420 scenarios, SHARP controlled FPR reliably (Fig. 6b), with CI lower
506 bounds below 5% in nearly all scenarios (Fig. 6c), matching the 5×2 paired t-test and the
507 empirical test of differences. Conclusions were consistent across datasets (Fig. S2), sample
508 sizes (Fig. S3), and an alternative simulation scheme (Fig. S4). A broader comparison across
509 additional tests (Methods Section 4.5, Supplementary Methods S8 & S9) confirmed the same
510 pattern: tests that effectively ignored fold dependence — a misapplication of the corrected
511 resampled t-test, the paired permutation test, and a bootstrap variant — had inflated FPR,
512 while tests that accounted for it — the 5×2 paired F-test and two other bootstrap variants —
513 controlled FPR reliably (Fig. S5).

514

515

516 **2.8 SHARP matches or exceeds the power of existing valid tests**

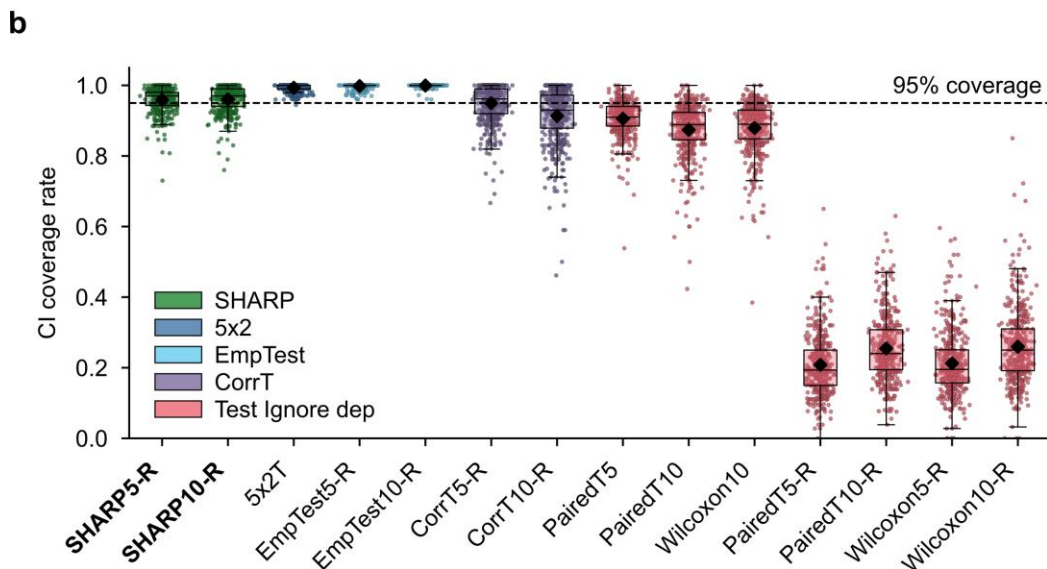
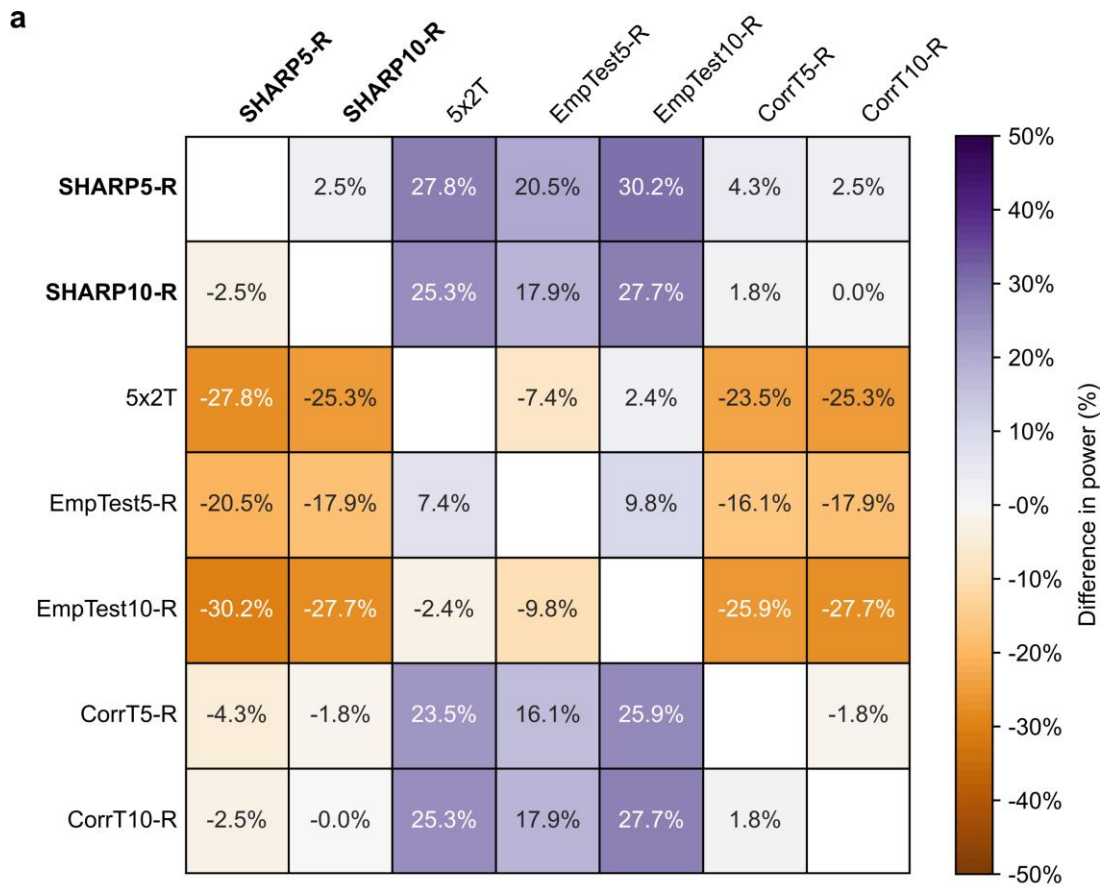
517 Aside from reliable FPR control, a statistical test should also have adequate power to detect
518 genuine performance differences. As argued in Section 2.6, SHARP may have higher power
519 than existing fold-aware tests. To test this prediction, we adapted the FPR simulation
520 procedure from Section 2.4.

521
522 Instead of corrupting both copies of the sampled dataset with the same level of noise (Fig.
523 5a), we corrupted only one copy and left the other clean (Fig. S6). The model trained on the
524 clean copy therefore had higher expected predictive performance by construction, and a
525 sensitive test should reject the null hypothesis of equivalent performance. As with FPR,
526 power was estimated within each of 420 scenarios as the fraction of repetitions in which the
527 null was rejected, then averaged across scenarios. More details are in Methods Section 4.4.

528
529 Fig. 7a shows the difference in average power between pairs of statistical tests. For example,
530 the SHARP test based on 5-fold cross-validation (SHARP5-R) outperformed the 5×2 paired
531 t-test by 27.8%, shown in the first row and third column of Fig. 7a. Overall, SHARP5-R had
532 the highest power, reflected in its entirely purple row in Fig. 7a. A wider set of tests is shown
533 in Fig. S7, where SHARP remained the most sensitive, followed by the corrected resampled
534 t-test. As in Section 2.7, we cross-checked our findings under an alternative simulation
535 scheme (Fig. S8; Supplementary Methods S6), which yielded the same conclusions, except
536 that SHARP5-R and CorrT5-R achieved similar power.

537
538 Together with Section 2.7, these results show that SHARP achieves substantially higher
539 power than the 5×2 paired t-test and the empirical test of differences, while matching them on
540 FPR control. Compared with the corrected resampled t-test, SHARP had similar or slightly
541 higher power, depending on the simulation scheme.

542



543
544

545 **Figure 7. Comparison of statistical power and confidence interval coverage across 420**
 546 **scenarios.** **a.** Difference in power between pairs of statistical tests, computed as the power of
 547 the statistical test (on the row) minus the power of the statistical test (on the column),
 548 averaged across 420 scenarios. A purple cell indicates that the test on the row achieved
 549 higher power than the test on the column; an orange cell indicates the opposite. For example,
 550 the cell in the first row and third column is 27.8% and purple, indicating that SHARP5-R has
 551 higher power than the 5×2 paired t-test. The tests shown here are a subset of the tests in Fig.

552 6, restricted to only valid statistical tests. SHARP5-R won every pairwise comparison, as
553 reflected in its entirely purple row. **b.** Comparison of 95% confidence interval (CI) coverage
554 rate across 420 scenarios. The CI coverage rate of a scenario is defined as the percentage of
555 non-overlapping subsamples in which the (estimated) true performance difference between
556 models fell inside the 95% CI of a given test. For a well-calibrated test, the coverage rate
557 should be exactly 95% (black dashed line). Each boxplot comprises 420 data points, each
558 representing the CI coverage rate of one scenario. The black dot indicates the mean coverage
559 rate across 420 scenarios. Tests that ignored fold dependence (red boxplots) had overly
560 narrow CIs, so their coverage rates were much lower than 95%.

561

562

563 **2.9 SHARP test demonstrates excellent confidence interval coverage**

564 Beyond FPR and statistical power, good statistical practice requires more than a binary
565 reject/fail-to-reject decision. A confidence interval communicates the magnitude and
566 precision of the estimated effect, essential for interpreting results in context. A statistically
567 significant difference may be too small to matter practically, while a non-significant result
568 with a wide interval may suggest insufficient power rather than true equivalence.

569

570 For a well-calibrated test, the 95% confidence interval should contain the true performance
571 difference between the two models 95% of the time (Methods Section 4.4.3). Across the 420
572 statistical power simulation scenarios, tests that ignored fold dependence produced overly
573 narrow CIs, resulting in coverage well below 95% (Fig. 7b). The empirical test of differences
574 and 5×2 paired t-test, which have low power (Fig. 7a), produced overly wide CIs with
575 coverage approaching 100%.

576

577 The corrected resampled t-test achieved a mean coverage of 95% under 5-fold cross-
578 validation but with wide variation across scenarios; this is consistent with its fixed between-
579 fold correlation that cannot adapt to the data (Section 2.6). SHARP5-R (96%) and SHARP10-
580 R (96%) both performed well. More tests are shown in Fig. S9. These conclusions also held
581 under an alternative simulation scheme (Fig. S10; Supplementary Methods S6), except that
582 the corrected resampled t-test now achieved stable coverage – the same scheme-dependence
583 noted for FPR in Section 2.7.

584

585 Across FPR, power, and confidence interval results, SHARP5-R best balances all three
586 criteria among the tests examined. The corrected resampled t-test is a close second on power,
587 but is less reliable on FPR control and CI coverage because it imposes assumptions about the
588 magnitude of fold-level correlations – assumptions that SHARP's design avoids.

589

590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630

3 Discussion

3.1 Practical guidance for valid model comparison

We have established that statistical tests that neglect cross-validation fold dependence are widely used in biomedical publications, and such tests inflate false positive rates. Here we consolidate these findings into practical guidance. Fig. 8a summarizes the statistical properties of the tests evaluated in this study, and Fig. 8b translates these properties into a decision procedure extending from the single-dataset case — the focus of this paper — to the multi-dataset case.

When only a single dataset is available, SHARP is the recommended default (Fig. 8). It provided the best overall balance across the regimes we examined: reliable false-positive control, high statistical power, and well-calibrated confidence intervals. The corrected resampled t-test is also a reasonable choice, particularly when sample size makes SHARP's split-half design less attractive. Further discussion about SHARP and corrected t-test are found in Section 3.5.

An alternative to cross-validation is a single train-validation-test split. Because the training set is fixed and the test set is used only once, test samples are independent and paired tests applied across test samples are valid. However, results can be highly sensitive to the particular split, (Varoquaux et al., 2017), so we do not recommend this procedure for small sample sizes.

A test within a single dataset supports inference only about the population from which the dataset was drawn. Access to multiple datasets enables broader inference. When 2–5 datasets are available, running a valid within-dataset test (SHARP or corrected t-test) on each dataset and checking for qualitative consistency provides *informal* evidence of cross-population superiority: a model that outperforms another across most or all datasets is more likely to be superior in new datasets drawn from the same populations (Fig. 8b).

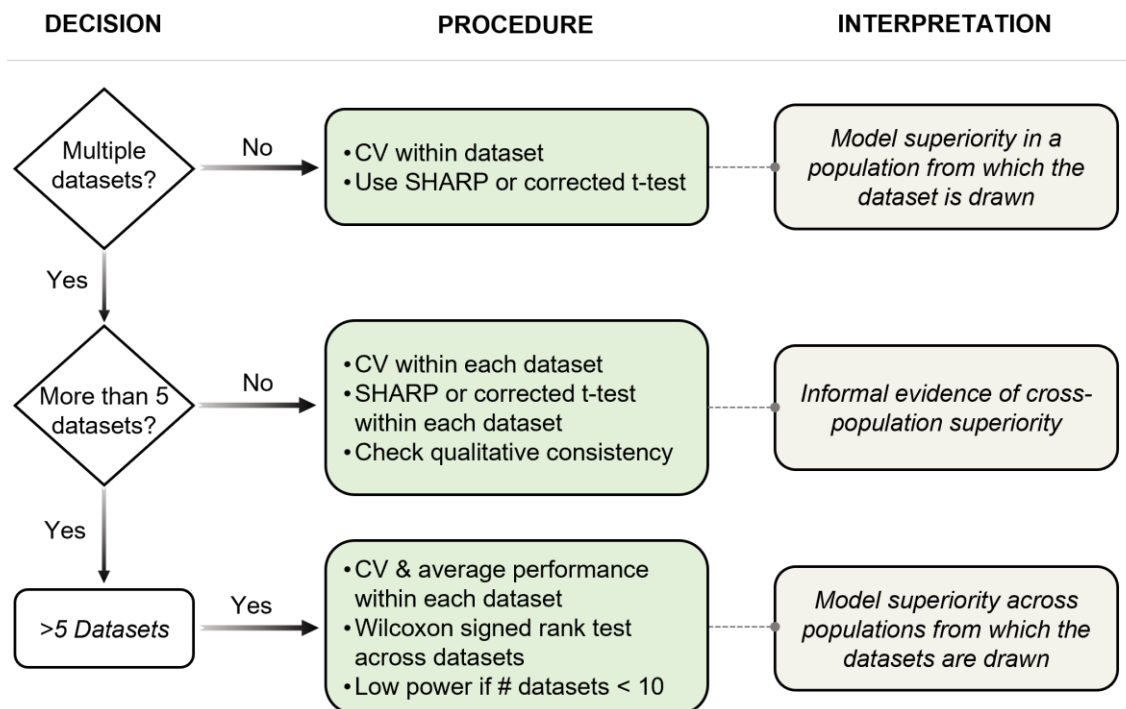
When more than five datasets are available, a Wilcoxon signed-rank test across dataset-level performance differences (Demšar, 2006) supports formal cross-population inference (Fig. 8b). However, power is low when there are fewer than ten datasets. Furthermore, access to six or more comparable datasets is often infeasible in many areas of biomedical research, where even two or three can be difficult to obtain.

It is also important to note that cross-population inference might often not be the goal (Dietterich, 1998). Many domain-specific studies aim to identify the best-generalizing model within a given dataset, for example, to deploy in the population from which that dataset was drawn. In such settings, a valid within-dataset test (SHARP or corrected t-test) is the appropriate tool, and the inference it supports – generalization to the source population – is exactly what the study requires.

a

Test	Fold-aware	FPR control	Power	CI coverage	Recommendation
Paired t-test	no	inflated	—	too narrow	not recommended
Wilcoxon signed-rank test	no	inflated	—	too narrow	not recommended within a dataset ¹
Permutation test (sign-flip)	no	inflated	—	too narrow	not recommended
DeLong's test	no	inflated	—	too narrow	not recommended
Bootstrap-t	no ²	inflated	—	too narrow	not recommended
Repetition-averaged corrected t-test	yes ³	inflated	—	too narrow	not recommended
Empirical test of differences	yes	controlled	lowest	too wide	valid but low power
Bootstrap-ET	yes	controlled	lowest	too wide	valid but low power
5×2 t-test	yes	controlled	lowest	too wide	valid but low power
5×2 F-test	yes	controlled	low	well-calibrated	valid but low power
Bootstrap-orig	yes	controlled	moderate	well-calibrated	valid but low power
Corrected paired t-test	yes	slightly inflated ⁴	high	slightly too narrow ⁵	reasonable alternative
SHARP	yes	controlled	high	well-calibrated	recommended default

b



632

633 **Figure 8. Guidelines for valid model comparison. a.** Statistical properties of tests for
634 comparing models in cross-validation. SHARP (highlighted) is recommended as the default
635 for within-dataset comparison; the corrected resampled t-test is a reasonable alternative,
636 especially when sample size is limited. **Footnotes:** ¹The Wilcoxon signed-rank test is invalid
637 within a dataset but valid for comparing dataset-level performance differences across datasets
638 (Demšar, 2006; see panel b). ²Bootstrap-t resamples fold-level differences under an implicit
639 independence assumption and is therefore invalid (see Methods Section 4.5.8). ³The
640 repetition-averaged corrected resampled t-test uses an overly liberal assumption of between-
641 fold correlation, leading to high FPR (see Supplementary Methods S9). ⁴FPR for the
642 corrected resampled t-test is controlled in some regimes but mildly elevated in others. ⁵95%
643 confidence interval for the corrected resampled t-test is well-calibrated in some regimes but
644 mildly elevated in others. Power is reported as "—" for tests that do not control false
645 positives at the nominal level, as power is not meaningfully defined in this case. **b.** Decision
646 framework for selecting an appropriate statistical comparison procedure based on dataset
647 availability. The flowchart guides users through sequential decisions (left) to arrive at a
648 recommended procedure (middle) and the corresponding scope of inference that the resulting
649 evidence supports (right).

650

651

652 **3.2 Invalid tests do not imply false findings**

653 The fold-dependence problem concerns uncertainty quantification, not point estimation.
654 Absent leakage or other circularity, cross-validated performance estimates remain valid
655 estimates of generalization to the population from which the dataset was sampled. What fold
656 dependence corrupts is the variance of those estimates — and hence p-values, confidence
657 intervals, and any test statistic that relies on assumed independence.

658

659 While our meta-analysis focused on 210 studies, many others were excluded due to lack of
660 methodological clarity, and these are unlikely to be free of invalid statistical tests. Although
661 we restricted our analysis to journals with impact factors of at least 15, there is little reason to
662 believe the issue is confined to this subset. Taken together, the scope of the fold-dependence
663 problem likely extends to thousands of studies.

664

665 Going forward, we encourage researchers to adopt the inference procedures recommended in
666 Fig. 8. However, for prior findings, an invalid statistical test does not imply a false positive:
667 sufficiently large effects may remain significant under more conservative procedures, and
668 any given model comparison may be only one component of a broader study whose overall
669 conclusions remain valid. We therefore do not claim that 97% of studies contain false
670 positives. Rather, 97% employed statistical tests that do not achieve the stated error control.
671 This distinction is critical, and prior findings warrant reassessment rather than wholesale
672 dismissal.

673

674 **3.3 Widespread invalid tests reflect a knowledge gap, not a rigor gap**

675 The dependence between cross-validation folds has been recognized for decades (Dietterich,
676 1998; Nadeau & Bengio, 2003; Bengio & Grandvalet, 2004), yet its neglect remains
677 pervasive. The absence of any association with impact factor, transparency policies, or open
678 science practices suggests a knowledge gap rather than a rigor gap. Existing safeguards —
679 code sharing, reporting checklists, and data availability — cannot detect violations of
680 statistical assumptions that are rarely taught and seldom reviewed. The issue is not that

681 researchers are cutting corners, but that the dependence structure induced by resampling is
682 not part of standard training in many biomedical disciplines.

683
684 Closing this gap will require changes at multiple levels. Graduate training should explicitly
685 cover the statistical properties of cross-validation. Peer review guidelines should flag cross-
686 validation-based statistical tests as requiring justification. Reporting standards should make it
687 straightforward for readers to assess whether valid inference has been performed (see next
688 section).

689
690 Intriguingly, researchers developing machine learning algorithms — as opposed to
691 biomedical researchers employing them — mostly report point estimates across multiple
692 datasets without formal statistical testing, despite established guidelines (Demšar, 2006; see
693 Fig. 8b). This practice raises a complementary challenge: when performance differences are
694 small, the absence of statistical inference makes it difficult to distinguish genuine
695 improvements from noise, as recently demonstrated in medical imaging artificial intelligence
696 (Christodoulou et al., 2025).

697 **3.4 Recommendations for reporting standards**

699 Our meta-analysis revealed substantial reporting deficiencies in the biomedical machine
700 learning literature beyond the issue of fold dependence. Of 1,875 studies that compared
701 model performance, 54% did not clearly apply a statistical test or report confidence intervals
702 — either because no formal inference was conducted, or because reporting was too vague to
703 tell. Among the remaining studies, many applied tests whose procedures were insufficiently
704 described to evaluate, so they also had to be excluded.

705
706 Studies employing cross-validation for model comparison should clearly report: (1) the
707 number of folds (e.g., 5-fold cross-validation), (2) the number of repetitions (e.g., 100
708 repeats), (3) the specific statistical test used, including whether it is one- or two-sided, and (4)
709 how fold-level statistics are aggregated, including the size of the resulting input to the
710 statistical test.

711
712 Clarifying the unit of analysis is particularly critical. For example, if 5-fold cross-validation
713 is repeated 100 times, are the performance differences averaged within each repetition,
714 yielding a vector of length 100 as input to the statistical test, or are all fold-level differences
715 concatenated, giving a vector of length 500? These choices are not interchangeable:
716 averaging within each repetition inflates the false positive rate of the misapplication of the
717 corrected resampled t-test (Fig. S5; Supplementary Methods S9).

718
719 Bootstrap methods warrant similar attention. In our meta-analysis, bootstrap procedures were
720 often described without sufficient detail to evaluate their validity. Authors should specify
721 whether bootstrapping is applied to generate training and test sets or to resample cross-
722 validation outputs, and how the resulting samples are used to compute p-values or confidence
723 intervals. Treating bootstrapped samples as independent observations massively inflates the
724 false positive rate, whereas using them as an empirical distribution yields valid inference
725 (Fig. S5) — albeit with low power (Fig. S7).

726
727 Confidence interval reporting requires the same care. In our meta-analysis, we excluded all
728 studies reporting only confidence intervals because too few specified how the intervals were
729 computed. This ambiguity is consequential: an interval based on the standard error of fold-
730 level statistics would incorrectly assume fold independence and be too narrow, while one

731 based on the 2.5th and 97.5th percentiles of those statistics would be overly conservative
732 (similar to the empirical test of differences). Authors should therefore specify the underlying
733 variance estimator or resampling scheme.

734

735 **3.5 SHARP and the corrected resampled t-test: practical considerations**

736 SHARP builds directly on standard cross-validation by introducing a split-half step that
737 eliminates overlap between paired statistics. In contrast to existing fold-aware tests, SHARP
738 estimates the relevant variance and correlation from the data rather than imposing strong
739 assumptions. We provide an open-source implementation compatible with scikit-learn,
740 enabling adoption with minimal changes to existing workflows.

741

742 The primary limitation of SHARP is the reduction in effective sample size induced by the
743 split-half procedure, which can be consequential when data are scarce. Different algorithms
744 or biomarkers can exhibit distinct scaling behavior, such that their relative performance
745 changes as training data increases. Because SHARP evaluates models using half the available
746 data, it may not fully reflect their relative performance at the full sample size.

747

748 The corrected resampled t-test is a reasonable alternative. Its FPR control is mildly inflated in
749 one simulation scheme (Fig. 6) but adequate in another (Fig. S4); its power is slightly below
750 SHARP in the first scheme (Fig. 7a) and similar to SHARP in the second (Fig. S8); and its CI
751 coverage is well-calibrated on average, though it shows substantial scenario-to-scenario
752 variability in one scheme (Fig. 7b), likely reflecting its assumption about the magnitude of
753 between-fold correlations.

754

755 Because the corrected resampled t-test uses the same inputs as standard cross-validation, it
756 can be adopted with no methodological overhead. We therefore recommend it as a reasonable
757 alternative to SHARP, especially when sample size is small.

758

759 **3.6 Broader implications for reproducibility**

760 A reproducibility crisis has been documented across multiple scientific fields, including
761 psychology (Open Science Collaboration, 2015), cancer biology (Errington et al., 2021), and
762 machine learning (Kapoor & Narayanan, 2023; Christodoulou et al., 2024). Common causes
763 include p-hacking, selective reporting, and underpowered studies. The mechanism described
764 here is different. Inflated false positive rates arise from a structural property of cross-
765 validation – correlated folds – rather than from researcher incentives or analytical flexibility.
766 Yet the consequence is similar: systematic overestimation of evidence that may fail to
767 replicate.

768

769 A distinguishing feature of this class of problems is that it is unusually tractable. Addressing
770 p-hacking or publication bias requires changes to incentives and research culture. By
771 contrast, correcting for fold dependence requires only adopting appropriate statistical tests or
772 evaluation procedures, which can be implemented without altering data collection or study
773 design.

774

775

776

4 Methods

777

778 4.1 Primer on cross-validation

779 In this section, we provide an overview of different variants of cross-validation. In K-fold
780 cross-validation, a dataset is divided into K non-overlapping partitions of data samples
781 (folds). For each of K iterations, one fold is treated as the test set, while training is performed
782 on the remaining folds. The procedure is repeated K times, so each fold has its turn to be the
783 test. When comparing two models, the entire K-fold cross-validation procedure can be
784 applied to each model separately. The K partitions are set up to be the same across the two
785 runs of K-fold cross-validation, so that there is fold-level correspondence between the two
786 models.

787

788 Within each fold, the performance of each model can be averaged across samples, yielding K
789 pairs of performance metrics for the two models. For each pair of performance metrics, the
790 difference can be computed, resulting in a vector \mathbf{D} of J performance differences, where $J =$
791 K . We note that the J performance differences are not independent because training sets
792 overlap across iterations (for $K > 2$) and the test set from one iteration contributes to the
793 training set of all other iterations. Many studies apply a statistical test to the vector \mathbf{D} of J
794 performance differences to evaluate the null hypothesis of equivalent prediction performance
795 between the two models.

796

797 For stability (Bouckaert & Frank, 2004; Varoquaux et al., 2017), many studies also repeat K-
798 fold cross-validation R times. In this scenario, we still obtain a vector \mathbf{D} of J performance
799 differences, but $J = K \times R$. We note that in this situation, there are two different fold-level
800 correlations here. First, the fold-level statistics are correlated within a single instance of K-
801 fold cross-validation. Across different repeats of K-fold cross-validation, there is now overlap
802 in both training and test sets, so fold-level statistics from different repeats of K-fold cross-
803 validation are correlated, and this correlation is likely different from the between-fold
804 correlation within a particular instance of K-fold cross-validation. Once again, many studies
805 apply a statistical test to the vector \mathbf{D} of J performance differences to evaluate the null
806 hypothesis of equivalent performance between the two models.

807

808 In Monte Carlo cross-validation, the datapoints that make up a dataset are repeatedly divided
809 into training and test sets randomly. If the performance metric is averaged across samples in
810 the test set and this procedure is repeated J times, we obtain a vector \mathbf{D} of J performance
811 differences. In this setup, there is overlap between both training and test sets across the J
812 repeats, so the fold-level statistics are non-independent. Similar to the previous setups, many
813 studies apply a statistical test to the vector \mathbf{D} of J performance differences to evaluate the null
814 hypothesis of equivalent performance between the two models.

815

816 In the above cross-validation schemes, the performance difference metric is averaged across
817 all samples in a test set, so the statistical test utilizes fold-averaged statistics. During our
818 meta-analysis, we observed some studies using "sample-level" instead of "fold-averaged"
819 statistics. To illustrate the distinction, consider a study performing a single instance of 10-
820 fold cross-validation on a dataset of 1,000 independent samples. A fold-averaged approach
821 averages performance differences within each fold, yielding 10 performance differences
822 metrics for the statistical test, i.e., the vector \mathbf{D} has 10 values. A sample-level approach
823 forgoes this aggregation and instead enters all 1,000 performance differences directly into the
824 statistical test, i.e., the vector \mathbf{D} has 1000 values. Correlations among performance
825 differences in the same fold might be stronger than correlation between folds (since the

826 trained model is the same within a fold), so sample-level statistics might further inflate the
827 false positive rates.

828
829 A final variant collapses the fold-level estimates in the opposite direction: rather than
830 disaggregating to individual samples, performance differences are averaged across folds
831 within each repetition of the cross-validation, so that each repetition contributes a single
832 value to the statistical test. We refer to these as “repetition-averaged” statistics, in parallel to
833 the “fold-averaged” and “sample-level” statistics above. For example, if 10-fold cross-
834 validation is repeated 30 times, a fold-averaged approach yields a vector D of length 300,
835 whereas a repetition-averaged approach first averages the 10 fold-averaged performance
836 differences within each repeat and then enters the resulting 30 values into the statistical test,
837 i.e., the vector D has 30 values. The correlation among the 30 repetition-averaged statistics is
838 larger than the correlation among the 300 fold-averaged statistics. Feeding this repetition-
839 averaged vector into the corrected resampled t-test (Nadeau & Bengio, 2003; Methods
840 Section 4.5.5) – which is calibrated for between-fold correlation, not the larger between-
841 repetition correlation – leads to a high false positive rate (Supplementary Methods S9).

842 843 **4.2 Meta-analysis**

844 4.2.1 Journal selection and search strategy

845 We conducted a literature search in PubMed on 18 September 2025 to identify original
846 research articles that compared machine learning models using cross-validation. The search
847 was restricted to articles published between 1 June 2020 and 1 June 2025. Eligible journals
848 were defined using Journal Citation Reports (JCR) and were required to have a Journal
849 Impact Factor (JIF) ≥ 15 in either 2023 or 2024. Review journals were excluded.

850
851 For each candidate journal, ISSN and eISSN information was obtained from JCR and
852 matched to the corresponding National Library of Medicine (NLM) journal abbreviation in
853 the NLM Catalog. Journals not indexed in the NLM Catalog were excluded, because PubMed
854 searches by journal required this identifier. For each eligible journal, the following PubMed
855 query was then executed:

856
857 ("outperform" OR "outperforms" OR "state-of-the-art" OR "benchmark" OR "benchmarking"
858 OR "algorithm comparison" OR "model comparison" OR "performance evaluation" OR
859 "method comparison") AND ("cross validation" OR "cross-validation" OR "crossvalidation")
860 AND ("2020/06/01"[Publication Date] : "2025/06/01"[Publication Date]) AND (Journal
861 Abbreviation [TA])

862
863 This search yielded 2,400 studies for subsequent screening (Fig. 2a).

864 865 4.2.2 Eligibility criteria and PRISMA workflow

866 Screening of the 2400 studies followed PRISMA guidelines (Page et al., 2021). A final set of
867 210 studies met the following four criteria.

868 869 1. **Direct model performance comparison.**

870 To satisfy this criterion, the study had to compare at least two distinct models – that
871 might differ in terms of algorithms, pipelines, hyperparameter configurations, feature
872 sets (e.g., biomarkers), or model-selection procedures – on the same prediction target.
873 Eligible studies must compare models based on “direct” performance metrics (e.g.,
874 accuracy, AUC, MSE, MAE, correlation, etc.) computed from comparing predictions
875 against ground truth. The models must be evaluated on the same test data. This

876 criterion resulted in 1875 papers involving model comparisons and 525 without model
877 comparisons (Fig. 2a).

878 We excluded studies that compared only downstream quantities derived from
879 predictions, as opposed to direct performance metrics. For example, a study might
880 develop two brain-age models from MRI and compare them based on which model's
881 brain age predictions better differentiate Alzheimer's disease from healthy
882 individuals, as opposed to comparing models' accuracy in predicting the training
883 target (chronological age). Although the fold dependence problem would also arise in
884 such cases, we excluded these studies because it was difficult to have a clear set of
885 criteria for what downstream analyses to include or exclude in the meta-analysis.
886

887 **2. Use of statistical inference for the comparison.**

888 The study had to report at least one p-value or confidence interval for a quantitative
889 comparison of machine learning models. Studies that only reported descriptive
890 statistics or performed statistical tests unrelated to model comparison (e.g., correlation
891 between model predictions and confounding variables) were considered ineligible.
892 This criterion resulted in 867 studies using statistical inference for comparing models,
893 while 1008 studies were eliminated.
894

895 **3. Fold-based statistical inference.**

896 Among the 867 studies, some did not use cross-validation, and a large fraction did not
897 clearly describe their inference procedure. To be included in our final analyses, a
898 study had to be clear that the values used to compute a statistical test or confidence
899 interval were derived from the test folds of cross-validation. For example, we
900 included studies that (i) performed K -fold cross-validation and used K pairs of values
901 for statistical inference; or (ii) conducted K repeated random train–test splits (i.e.,
902 Monte Carlo cross-validation) and used K pairs of values for statistical inference.
903 There were also studies that repeated K -fold cross-validation multiple times. For
904 example, a study that repeated 10-fold cross-validation 100 times and used 1,000
905 pairs of values for a t-test, would be included in the meta-analysis. If such studies
906 averaged performance across the 10 folds within each repetition and used the
907 resulting 100 paired values in a t-test, they would also be included. The studies also
908 had to clearly state the name of the statistical test. This criterion resulted in a set of
909 265 studies.
910

911 **4. Use of statistical tests that ignore or account for fold dependence.**

912 Of the 265 studies, one employed a permutation test that was invalid for reasons
913 unrelated to fold dependence; we therefore considered it outside the scope of the
914 current study (see Supplementary Results for details). We also excluded studies that
915 reported only confidence intervals, as most simply asserted the 95% interval without
916 specifying how it was computed, a problem that has been independently documented
917 in medical imaging (André et al., 2026). A 95% confidence interval based on the
918 standard error of fold-level statistics would incorrectly assume fold independence,
919 yielding an interval that is too narrow. Conversely, a 95% interval computed from the
920 2.5th and 97.5th percentiles of fold-level statistics approximates the empirical test of
921 differences (Methods Section 4.5.7) and would be overly conservative. Without
922 further methodological details, these studies could not be reliably analyzed. Applying
923 these criteria resulted in a final set of 210 studies.

924

925 The PRISMA flow diagram describing identification, screening, eligibility assessment and
926 inclusion is shown in Fig. 2a.

927

928 4.2.3 Scientific field classification

929 To explore the prevalence of invalid statistical inference across scientific fields, we assigned
930 each study to a scientific field. Scientific fields were defined using the Web of Science
931 Journal Citation Reports (JCR) subject categories. JCR listed 254 categories, and each
932 journal could belong to one or more categories. Because our PubMed-based search targeted
933 biomedical and life-science content, not all JCR categories were represented. For journals
934 assigned to one scientific category, all studies published in the journal were assigned to that
935 category. For studies published in journals assigned to multiple categories (e.g., Nature), we
936 used a large language model (LLM) to determine the most relevant category for each study.
937 See Section 4.2.4 for more details about the LLM procedure, as well as the manual checks
938 performed by human raters.

939

940 4.2.4 LLM-assisted screening and manual human verification

941 To scale screening and data extraction, we used an LLM (Claude Opus 4.1, Anthropic) as an
942 assistant in combination with manual human verifications. For each of 2400 studies, we
943 provided the main text scraped from PubMed, including figure captions, but without figures,
944 tables or supplementary material, together with a structured prompt. The full prompt text is
945 provided in Supplementary Methods S1.1. The LLM was asked to answer five questions for
946 each study:

947

- 948 1. **Scientific field classification (Q1).** If a study was published in a multi-field journal, the
949 LLM was instructed to assign the study to a single scientific subfield using a predefined
950 set of categories, or “OTHERS” if none applied.

951

952 *Manual verification:* In a random audit of 50 studies, two human raters (TZ and HL)
953 manually checked the studies. Any discrepancy was resolved by discussion and
954 consensus. The LLM was in 100% agreement with the human raters. Therefore, we used
955 the LLM results directly.

956

- 957 2. **Model comparison evaluation (Q2).** The LLM was asked to determine whether the
958 study satisfied PRISMA criterion 1 (Methods Section 4.2.2).

959

960 *Manual verification:* In a random audit of 50 studies, two human raters (TZ and HL)
961 manually checked the studies. Any discrepancy was resolved by discussion and
962 consensus. The LLM was in 100% agreement with the human raters. Therefore, we used
963 the LLM results directly.

964

- 965 3. **Statistical inference (Q3).** If the LLM’s answer to Q2 was “Yes”, then the LLM was
966 instructed to also identify every instance where a p-value or confidence interval was
967 reported when comparing model performance. Articles with at least one instance were
968 deemed eligible under PRISMA criterion 2 (Methods Section 4.2.2).

969

970 *Manual verification:* In a random audit of 50 studies, two human raters (TZ and HL)
971 manually checked the studies. Any discrepancy was resolved by discussion and
972 consensus. The LLM agreed with two human raters (TZ and HL) on 39 studies (78%).
973 For the remaining 11 studies, there was only one false negative, i.e., LLM counted a

974 study as being ineligible under PRISMA criterion 2, but the human raters disagreed. We
975 considered this level of false negatives to be acceptable. The remaining 10 errors were
976 false positives, i.e., the LLM counted the studies as being eligible under PRISMA
977 criterion 2, but human raters disagreed. Given the relatively high false positive rate,
978 human raters read all 1197 studies that the LLM considered to satisfy PRISMA criterion
979 2. More specifically, each study was manually examined by two human raters (TZ, HL or
980 SZ), and any discrepancy was resolved by discussion and consensus. After going through
981 all studies, we ended up with 867 studies that satisfied PRISMA criterion 2.
982

- 983 4. **Statistical test details (Q4).** For each instance identified in Q3, the LLM was instructed
984 to extract the name of the test or confidence interval, details about compared models,
985 performance metric(s) used for the comparison, and the cross-validation procedure used.
986 The LLM was asked to support its conclusions with direct quotations from the study.
987

988 *Manual verification:* The extracted quotations served as a reference for human
989 verification. However, we note that every study was manually examined by two human
990 raters (TZ, HL or SZ) to re-derive the above information. Any discrepancy between
991 human raters was resolved by discussion and consensus.
992

- 993 5. **Data classification (Q5).** For each instance identified in Q3, the LLM was also asked to
994 classify whether the values entering the test were derived from cross-validation folds
995 (referred to as “resampling units” in our prompt) or as “Other/Unclear”, corresponding to
996 PRISMA criterion 3 (Methods Section 4.2.2).
997

998 *Manual verification:* The LLM classification served as a reference for human
999 verification. However, we note that every study was manually examined by two human
1000 raters (TZ, HL or SZ) to determine whether PRISMA criterion 3 was satisfied. Any
1001 discrepancy between human raters was resolved by discussion and consensus.
1002

1003 For each of the final set of 210 studies that satisfied all four PRISMA criteria, the human-
1004 verified information from Q4 and Q5 above was used to classify the article into one of two
1005 groups: (i) statistical tests ignoring fold dependence, and (ii) statistical tests accounting for
1006 fold dependence.
1007

1008 4.2.5 Journal policy rigor scoring

1009 To evaluate whether journal guidelines might influence the prevalence of invalid statistical
1010 inference, we defined a journal-level “policy rigor” score comprising two domains with a few
1011 criteria under each domain:
1012

1013 1. **Transparency and reproducibility**

- 1014 ○ *Code availability policy:* requirement or encouragement to share analysis code
1015 or model scripts.
- 1016 ○ *Data availability policy:* requirement or encouragement to share underlying or
1017 processed data.
- 1018 ○ *Data/code availability statement requirement:* requirement or encouragement
1019 for a formal data and/or code availability statement in the manuscript.
- 1020 ○ *Code for peer review:* requirement or encouragement to provide custom code
1021 to editors and reviewers during peer review.

1022 2. **Statistical rigor and methodological guidance**

- 1023 ○ *Statistical test guidance*: explicit guidance on appropriate statistical testing
1024 (for example, paired vs unpaired tests, reporting of uncertainty, comparison
1025 under cross-validation).
- 1026 ○ *Statistical review*: whether dedicated statistical review is part of the editorial
1027 or peer-review process.
- 1028 ○ *Reporting standards or checklists*: requirement or encouragement to follow
1029 established methodological reporting guidelines (for example, TRIPOD-AI,
1030 CONSORT-AI, PRISMA).

1031
1032 In the case of transparency and reproducibility, journals were scored on a three-point scale
1033 for each of the four criteria: 0 = not mentioned or no guidance; 1 = encouraged or optional; 2
1034 = mandatory or strongly enforced. Since there were four items, the total score ranged from 0
1035 to 8.

1036
1037 In the case of statistical rigor and methodological guidance, journals were also scored on a
1038 three-point scale for each criterion with 0 corresponding to least rigor and 2 corresponding
1039 with most rigor. For more details about the scoring criteria, refer to Supplementary Table S1.
1040 Since there were three items, the total score ranged from 0 to 6.

1041
1042 To derive the scores, we used an LLM (GPT-5, OpenAI) to search each journal's official
1043 webpages (Instructions for Authors, Editorial Policies, Submission Guidelines, Reporting
1044 Standards and related pages) on November 7, 2025 to identify the most authoritative policy
1045 statements. We changed the LLM from Anthropic to OpenAI, because the OpenAI LLM
1046 interface was more suitable for open web search at the time of the search. The LLM prompt
1047 is provided in Supplementary Methods S1.2.

1048
1049 Scoring from the LLM were verified by two human raters (TZ and HL). Any discrepancy
1050 between human raters was resolved by discussion and consensus. Finally, scores were
1051 summed across the two domains, resulting in an overall journal rigor score that ranged from 0
1052 to 14. The score was used in analyses presented in Fig. 4. See Methods Section 4.2.7 for
1053 details about the statistical analyses.

1054 4.2.6 Study-level code and data availability

1055 For each study, we used an LLM (Claude Opus 4.1, Anthropic) to determine whether the
1056 study provided data or code. Similar to Methods Section 4.2.4, for each of 210 studies, we
1057 provided the main text scraped from PubMed, including figure captions, but without figures,
1058 tables or supplementary material, together with a structured prompt. The full prompt text is
1059 provided in Supplementary Methods S1.3. In a random audit of 50 studies, the LLM was in
1060 100% agreement with two human raters (TZ and HL). Therefore, the LLM results were used
1061 in analyses presented in Fig. 4. See Methods Section 4.2.7 for details about the statistical
1062 analyses.

1063 4.2.7 Test for temporal trend & association with scientific field, impact factor & scientific 1064 rigor

1065 We tested for a trend in the publication rate of the 210 studies over time. Letting Y_i be
1066 number of studies published in the i -th 6-month period, we fit a log-linear Poisson regression
1067 as follows:

$$1070 Y_i \sim \text{Poisson}(\mu_i), \quad \ln \mu_i = \beta_0 + \beta_1 \times i; \quad i = 0, \dots, 9,$$

1071

1073 where β_1 quantifies the trend over time and β_0 is the intercept. If $\beta_1 = 0$, then this would
1074 suggest that the number of studies stayed constant over time. If β_1 is positive, then this means
1075 that volume of relevant studies is increasing. For example, in the current study, we found that
1076 $\beta_1 = 0.07$ ($p = 6e-3$), so $e^{0.07} = 1.070$. This suggests that the volume of relevant studies
1077 increased by approximately 7.0% per 6-month period, or 14.4% per year. The Poisson
1078 regression model was fitted using python package statsmodels 0.14.5 with a p-value based on
1079 a Wald test.

1080

1081 We tested for the association between the study attribute “ignored fold dependence” and
1082 various factors. For ordinal factors of time, impact factor and journal rigor, we used a
1083 permutation-based Cochran–Armitage trend test. We did not use logistic regression for these
1084 analyses because the proportions of studies ignoring fold dependence were close to one,
1085 leading to a high prevalence of zero- and low-rate cells. The Cochran–Armitage test is
1086 sensitive to monotonic differences in the rate of studies that ignored fold dependence.
1087 Specifically, we constructed a $N_{rows} \times 2$ contingency table in which rows indexed ordered
1088 strata (time intervals, impact-factor bins, or rigor levels) and columns counted studies that
1089 ignored versus accounted for fold dependence. Under the null hypothesis of no monotonic
1090 trend, we generated an empirical null distribution by repeatedly sampling 100,000
1091 contingency tables conditional on the observed marginal totals (fixed row and column sums).
1092 For each sampled table we computed the Cochran–Armitage trend statistic. The two-sided p-
1093 value was estimated as the proportion of sampled Cochran-Armitage trend statistics whose
1094 absolute value was greater than or equal to the absolute value of the observed statistic.

1095

1096 For the nominal factors of scientific field, study-level code availability and study-level data
1097 availability, we used a permutation-based analogue of the Fisher–Freeman–Halton exact test.
1098 We first constructed a $N_{rows} \times 2$ contingency table with rows representing levels of the
1099 nominal factor and the same two outcome columns as above. Under the null hypothesis of no
1100 difference between fields, we generated 100,000 tables conditional on the observed margins.
1101 For each sampled table, we computed the Pearson χ^2 statistic as a measure of deviation from
1102 independence, forming an empirical null distribution. The p-value was estimated as the
1103 proportion of permuted χ^2 values greater than or equal to the observed χ^2 .

1104

1105 **4.3 Datasets**

1106 We evaluated the performance of statistical tests and confidence interval estimation using
1107 four datasets: EMNIST Digits (Cohen et al., 2017), UK Biobank (UKB; Alfaro-Almagro et
1108 al., 2018), Covertypes (Blackard, 1998), and KEGG Metabolic Pathway (Muhammad Naem,
1109 2011). These datasets spanned image recognition, neuroimaging, ecological classification,
1110 and systems biology, enabling a comprehensive comparison of statistical inference
1111 approaches across distinct data modalities and learning tasks (classification and regression).

1112

1113 **4.3.1 EMNIST**

1114 We used the “digits” subset of the Extended MNIST (EMNIST) dataset (Cohen et al., 2017),
1115 which comprised 240,000 grayscale images (28×28 pixels) in the training set and 40,000
1116 grayscale images (28×28 pixels) in the test set. For the current study, we only utilized the
1117 training set. There were 24,000 images per digit class (0–9). The machine learning task was
1118 to classify each image into one of the 10-digit classes.

1119

1120 **4.3.2 UK Biobank**

1121 We analyzed data from 36,454 UK Biobank participants (Alfaro-Almagro et al., 2018),
1122 consistent with our previous studies (He et al., 2022; Wulan et al., 2024). The prediction task

1123 was regression, with age (defined as MRI scan date minus birth year and month) as the target
1124 variable. The volumes of 101 cortical and subcortical gray-matter regions generated by the
1125 FreeSurfer software (Fischl, 2012) were features used for prediction. Raw volumetric
1126 measures were first normalized by dividing by the intracranial volume of each participant.
1127 Feature z-normalization was then performed using the mean and standard deviation computed
1128 from a held-out subset of 454 randomly sampled participants, which were subsequently
1129 applied to the remaining 36,000 participants. This yielded a final dataset of 36,000
1130 individuals, each represented by 101 standardized features.

1131 4.3.3 Covertypes

1133 The Covertypes dataset enabled the benchmarking of forest cover type classification based on
1134 cartographic variables derived from remote sensing and US Forest Service data. Each
1135 instance corresponded to a 30×30 meter cell in one of four wilderness areas within the
1136 Roosevelt National Forest, Colorado (Blackard, 1998).

1137
1138 The dataset contained 581,012 instances and 54 features, including both continuous variables
1139 (elevation, slope, aspect, hillshade values, distance to hydrology/roads) and binary variables
1140 encoding soil types and wilderness area indicators. Labels represented one of seven forest
1141 cover types, derived from the USFS Region 2 Resource Information System (RIS).

1142
1143 Following previous studies (Collobert et al., 2001), we constructed a binary classification
1144 task: class 1 included all cover types except for type 2, while class 2 included only cover type
1145 2 (commonly associated with lodgepole pine in Rawah and Comanche Peak). This
1146 transformation resulted in a more balanced distribution: 297,711 samples in class 1 and
1147 283,301 samples in class 2.

1148 4.3.4 KEGG Metabolic Pathway

1149 The KEGG (Kyoto Encyclopedia of Genes and Genomes) Metabolic Pathway database
1150 describes interactions between biochemical compounds, enzymes, and genes. These
1151 pathways can be represented as graphs using two models: (i) the reaction network, where
1152 nodes correspond to substrates and products, and edges correspond to catalyzing genes or
1153 enzymes; and (ii) the relation network, where compounds form edges between gene/enzyme
1154 nodes (Muhammad Naeem, 2011). A large set of pathway entries were parsed and
1155 transformed into graphs using both representations.

1156
1157
1158 The dataset provided graph-based features extracted using Cytoscape (Shannon et al., 2003),
1159 a software platform for biological network visualization and analysis. There were a total of
1160 53,414 samples, each described by 20 topological features (e.g., degree centrality,
1161 betweenness centrality, and clustering coefficient). In the current study, the learning objective
1162 was to predict the clustering coefficient of each network instance, making this a regression
1163 task.

1164 4.4 Simulations to evaluate FPR, statistical power and confidence intervals

1166 4.4.1 Simulation procedure for evaluating false positive rates (FPR)

1167 We conducted analyses on four datasets — EMNIST, UKB, Covertypes, and KEGG
1168 Metabolic Pathway — using four sample sizes: $N = 100$, 500, 1,000, and 2,000. For
1169 EMNIST, $N = 100$ was excluded because with 10 digit classes, it was not possible to
1170 guarantee at least one instance of each class in the test set under our proposed SHARP test,
1171 which required a split-half step within the cross-validation scheme (see Methods Section 4.6).
1172 For UKB, $N = 2,000$ was excluded because the maximum number of non-overlapping

1173 sampled datasets that could be drawn from UKB at this sample size was fewer than 20 (i.e.,
1174 $36,000 / 2,000 = 18$), which we considered too few to yield stable false positive rate
1175 estimates.

1176
1177 We drew m non-overlapping sampled datasets of size N from the full dataset. The value of m
1178 was set to 100 when feasible; otherwise, it was set to the maximum number of non-
1179 overlapping sampled datasets that could be drawn. For example, for the UK Biobank dataset
1180 at $N = 1,000$, this yielded $m = 36$ (i.e., $36,000 / 1,000$).

1181
1182 For each sampled dataset of size N , two noisy versions were generated by independently
1183 injecting the same level of random noise. Noise was introduced by permuting the target
1184 variable (e.g., digit labels for EMNIST) for a fixed proportion of samples. These were
1185 referred to as Noisy Dataset 1 and Noisy Dataset 2. We tested 10 permutation percentages in
1186 total: 10% to 100% with an increment of 10%.

1187
1188 For each noisy version of the sampled dataset of size N , we applied (i) a single instance of K -
1189 fold cross-validation or (ii) multiple repeats of K -fold cross-validation or (iii) Monte Carlo
1190 cross-validation. In classification tasks, stratified splitting was used to preserve label
1191 proportions across folds. All data splits were kept identical across the original sampled
1192 dataset and the two noisy datasets to ensure correspondence for paired statistical tests. We
1193 note that the SHARP test utilized a different cross-validation scheme (see Methods Section
1194 4.6 for details).

1195
1196 A single algorithm was selected for each dataset based on initial exploration using the
1197 PyCaret Python package (Ali, 2020): logistic regression classifier for EMNIST, ridge
1198 regression for UK Biobank, Extra Trees classifier for Covtype, and Extra Trees regression
1199 for KEGG Metabolic Pathway (see Methods Section 4.4.4). For each noisy version of the
1200 sampled dataset, the algorithm was trained on the training folds and evaluated on the test fold
1201 of the original (unpermuted) sampled dataset. Evaluating both models on the same
1202 unpermuted test data ensured fair comparison between the two noisy-trained models.

1203
1204 For every training-test split, a pair of fold-level performance metrics were computed for each
1205 pair of noisy datasets, and the model performance differences between the two noisy datasets
1206 were computed, resulting in a vector \mathbf{D} comprising J fold-level differences (except for
1207 SHARP; see Methods Section 4.6). For example, if K -fold cross-validation was repeated R
1208 times, then vector \mathbf{D} will be of length $J = K \times R$. For a given two-sided statistical test (e.g.,
1209 paired t-test), the vector \mathbf{D} was used to evaluate the null hypothesis $\mu = 0$, where μ is the true
1210 expected difference between the pair of trained models.

1211
1212 Since the noise level and algorithm were held constant across the pair of noisy datasets, the
1213 two trained models had identical expected predictive performance. Any rejection of the null
1214 hypothesis constituted a false positive. The false positive rate (FPR) was estimated as the
1215 fraction of m sampled datasets in which the null hypothesis was rejected at a significance
1216 threshold of 0.05. A well-calibrated test should reject the null hypothesis no more than 5% of
1217 the time. Recall that the m datasets were non-overlapping, so traditional binomial methods
1218 can be used to compute confidence intervals for the FPR. Specifically, the 95% confidence
1219 interval for the FPR was computed using the Wilson score interval with continuity correction
1220 (Newcombe, 1998); the formula is provided in Supplementary Methods S2.

1221

1222 Depending on the analysis, K -fold cross-validation was either performed once ($R = 1$) or
1223 repeated R times. When repeated, we targeted a total of 300 folds ($K \times R = 300$) to ensure
1224 convergence of the statistical tests — for example, $R = 30$ for 10-fold and $R = 60$ for 5-fold
1225 cross-validation, so J was 300. For Monte Carlo cross-validation, we always used an 80/20
1226 train-test split repeated 300 times, so J was again 300. The exception was the 5×2 statistical
1227 tests, which by design required 2-fold cross-validation repeated 5 times, so $J = 10$.

1228
1229 We note that given 10 noise levels (10% to 100%) and sample size ($N = 100$ to 2,000) across
1230 four datasets, there were 140 conditions. For each condition, we considered three
1231 hyperparameter-selection strategies: (1) fixed hyperparameters, (2) hyperparameter selection
1232 via a single training-validation split of the training data, and (3) nested cross-validation, in
1233 which hyperparameters were selected using cross-validation within the training data (see
1234 Methods Section 4.4.4). This resulted in $140 \times 3 = 420$ simulation scenarios in total. The FPR
1235 results are reported in Figs 5 and 6.

1236

1237 4.4.2 Simulation procedure for evaluating statistical power

1238 The simulations to benchmark statistical power followed the same procedure as those for
1239 FPR (Section 4.4.1), with one key difference. For FPR, the same algorithm was trained on
1240 two equivalently noisy datasets, so that the resulting pair of models was expected to have
1241 equivalent predictive performance. To assess statistical power, the same algorithm was
1242 instead trained on one noisy dataset and one original (unpermuted) sampled dataset. In this
1243 case, the model trained on the original (clean) sampled dataset was expected to outperform
1244 the model trained on the noisy dataset. A sensitive test should reject the null hypothesis
1245 $H_0: \mu = 0$ in a two-tailed test, where μ is the true expected difference in predictive
1246 performance between the two trained models.

1247

1248 As in the FPR case, power was estimated as the fraction of m datasets in which the null
1249 hypothesis was rejected at a significance threshold of 0.05. Crucially, a rejection was counted
1250 as a true positive only when the original-dataset model was found to significantly outperform
1251 the noisy-dataset model; cases in which the test rejected the null in the opposite direction
1252 were not counted as successes. Power was then averaged across 420 scenarios, yielding an
1253 average power for each statistical test.

1254

1255 4.4.3 Confidence Intervals for Performance Difference

1256 While there are different methods to compute a confidence interval (CI), here we used test-
1257 inversion which defines the CI by the set of null-hypothesis parameter values that cannot be
1258 rejected.

1259

1260 Recall that there were 420 simulation scenarios with each scenario corresponding to a given
1261 dataset (e.g., EMNIST), sample size (e.g., 2000), noise permutation percentage (e.g., 20%)
1262 and hyperparameter setting (e.g., fixed hyperparameters). For each scenario, there were m
1263 (non-overlapping) sampled datasets. For each sampled dataset, we injected noise to it,
1264 resulting in a noisy dataset. In the power simulations (Methods Section 4.4.2), the same
1265 algorithm was trained on one noisy dataset and one original (unpermuted) sampled dataset.
1266 The model trained on the original sampled dataset was expected to outperform the model
1267 trained on the noisy data.

1268

1269 Following the power simulations (Methods Section 4.4.2), for a given pair of original and
1270 noisy datasets, we performed cross-validation to obtain a vector \mathbf{D} of J fold-level prediction
1271 performance differences, with mean \bar{D} . We applied various statistical tests to the vector \mathbf{D} to

1272 evaluate the null hypothesis $\mu = 0$, where μ is the true expected difference between the pair
1273 of trained models.

1274

1275 To obtain the 95% confidence interval for μ for a given statistical test (e.g., paired t-test), we
1276 tested null hypotheses $\mu = \mu_0$ for different values of μ_0 to discover the range of values such
1277 that $p \geq 0.05$. This search was achieved by starting with $\mu_0 = \bar{D}$, and performing a binary
1278 search in the direction larger than \bar{D} and in the direction smaller than \bar{D} . The set of non-
1279 significant parameter values, $[\mu_L, \mu_U]$, is the 95% confidence interval for the performance
1280 difference between the two models.

1281

1282 Ideally, we would evaluate the accuracy of the computed 95% confidence intervals by
1283 comparing with the true difference μ_{true} , but this is not observable. Instead, we averaged \bar{D}
1284 across all m non-overlapping sets, resulting in the estimate $\hat{\mu}_{\text{true}}$. Across the m confidence
1285 intervals, we counted the fraction of times $\hat{\mu}_{\text{true}}$ fell within the confidence intervals, which
1286 we referred to as the overall coverage rate. A well-calibrated 95% confidence interval will
1287 cover μ_{true} 95% of the time. The whole procedure was repeated 420 times, once for each
1288 scenario.

1289

1290 4.4.4 Choice of algorithms, metrics and hyperparameter tuning schemes

1291 To determine an algorithm for each dataset, we used the python package PyCaret (Ali, 2020)
1292 to explore a set of algorithms available in the scikit-learn package with default
1293 hyperparameters (Pedregosa et al., 2011). We preferred algorithms that performed well on the
1294 original data, ensuring room for performance to decrease when trained on noisy data under
1295 our FPR and power simulation schemes. Computationally intensive algorithms with good
1296 performance were excluded (e.g., CatBoost).

1297

1298 We considered three hyperparameter tuning strategies. The first strategy used fixed
1299 hyperparameters, where a single predetermined configuration was applied based on scikit-
1300 learn defaults. The second strategy used an internal 80%/20% train-validation split within the
1301 training set. For each candidate hyperparameter value, the model was trained on the 80%
1302 subset and evaluated on the 20% validation set. The training-validation split was fixed across
1303 candidates, so validation performance was comparable. The hyperparameter value with the
1304 best validation performance was then used to train a final model on the full training set,
1305 which was applied to the test set.

1306

1307 The third strategy used nested cross-validation. For each candidate hyperparameter value, 2-
1308 fold cross-validation was performed on the training set. The cross-validation split was fixed
1309 across candidates, so performance was comparable. The hyperparameter value with the best
1310 cross-validation performance was then used to train a final model on the full training set,
1311 which was applied to the test set.

1312

1313 For the EMNIST dataset, where the goal was to classify images into one of 10-digit classes,
1314 we used multi-class logistic regression classifier. The performance metric was classification
1315 accuracy, defined as the fraction of samples in the test set that was classified correctly. The
1316 hyperparameter of interest was C , the inverse regularization strength. In the fixed
1317 hyperparameter setting, C was set to 1.0, while in the hyperparameter tuning settings, C was
1318 searched across the following values: [0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5,
1319 10, 50, 100, 500, 1000, 5000, 10000].

1320

1321 For the UKB dataset, where the goal was to predict age of the participant, we applied linear
1322 ridge regression. The performance metric was the coefficient of determination in the test set
1323 (Wright, 1921). The hyperparameter of interest was the regularization parameter α . In the
1324 fixed hyperparameter setting, α was set to 10, and for hyperparameter tuning, α was searched
1325 across the following values: [0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 1.5, 2, 3, 4,
1326 5, 10, 15, 20, 30, 40, 50, 100, 150, 200, 300, 400, 500, 1000, 2000, 5000, 10000].
1327

1328 For the Covertypes dataset, where the goal was to classify samples into one of two cover type
1329 classes, we used the Extra Trees classifier. The performance metric was classification
1330 accuracy, defined as the fraction of samples in the test set that was classified correctly. In the
1331 fixed hyperparameter setting, we used scikit-learn defaults with 100 estimators. For
1332 hyperparameter tuning, ranges were chosen following prior studies (Komer et al., 2014;
1333 Grinsztajn et al., 2022) and optimized using the Tree-structured Parzen Estimator (TPE;
1334 Bergstra et al., 2011, 2013).
1335

1336 For the KEGG Metabolic Pathway dataset, where the goal was to predict the clustering
1337 coefficient of each network sample, we applied the Extra Trees regressor. The performance
1338 metric was the mean squared error, computed by averaging the squared error across all test
1339 samples. In the fixed hyperparameter setting, we used scikit-learn defaults with 100
1340 estimators. For hyperparameter tuning, ranges were chosen following prior studies (Komer et
1341 al., 2014; Grinsztajn et al., 2022) and optimized using the Tree-structured Parzen Estimator
1342 (TPE; Bergstra et al., 2011, 2013).
1343

1344 Across the four datasets, we evaluated four sample sizes each, with the exception of EMNIST
1345 and UKB, for which only three sample sizes each were considered (see Methods Section
1346 4.4.1). Combined with ten noise levels and three hyperparameter optimization schemes, this
1347 yielded $(4 + 4 + 3 + 3) \times 10 \times 3 = 420$ simulation scenarios in total.
1348

1349 The KEGG Metabolic Pathway dataset was the only dataset in which performance was
1350 recorded at the individual sample level (squared error per test sample), rather than aggregated
1351 across the entire test set. This made it the only dataset suitable for comparing fold-averaged
1352 and sample-level statistical tests (Results Section 2.5). For this analysis, the combination of
1353 four sample sizes, ten noise levels, and three hyperparameter optimization schemes yielded 4
1354 $\times 10 \times 3 = 120$ scenarios.
1355

1356 **4.5 Existing hypothesis testing approaches for cross-validation**

1357 We consider the problem of testing whether one machine learning model statistically
1358 outperforms another machine learning model in a given dataset based on cross-validation. For
1359 example, in K -fold cross-validation, the dataset of N samples is partitioned into K mutually
1360 exclusive folds. Each fold serves once as a test set, while the remaining $K-1$ folds form the
1361 training set. The training set is used to train each model, and model performance is evaluated
1362 in the test set. This yields K pairs of performance metrics on which we want to perform a
1363 statistical test.
1364

1365 If we repeat K -fold cross-validation R times, then we have $J = K \times R$ pairs of performance
1366 metrics. In general, we assume the statistical test operates on a vector \mathbf{D} (of length J), where
1367 each element of vector \mathbf{D} corresponds to a fold-level performance difference between the two
1368 models for a particular test fold.
1369

1370 For example, in the case of performing 10-fold cross-validation once, the vector \mathbf{D} will be of
1371 length $J = 10$. If we perform 10-fold cross-validation 30 times, vector \mathbf{D} will be of length $J =$
1372 300. In the case of 80-20 Monte Carlo cross-validation where the dataset is repeatedly split
1373 into 80% training set and 20% test set 300 times, the vector \mathbf{D} will again be of length $J = 300$.

1374

1375 In the following subsections, we briefly summarize various statistical tests used in the
1376 literature. We note that the resampled paired t-test (Section 4.5.1), Wilcoxon signed-rank test
1377 (Section 4.5.2), permutation test (Section 4.5.3) and DeLong’s test (Section 4.5.4) are
1378 expected to have elevated FPR. The corrected resampled t-test (Section 4.5.5), the 5×2 tests
1379 (Section 4.5.6) and the empirical test of differences (Section 4.5.7) implicitly or explicitly try
1380 to account for fold dependence. The bootstrap may or may not be valid depending on how it
1381 was implemented (Section 4.5.8).

1382

1383 4.5.1 Resampled paired t-test

1384 As demonstrated in our meta-analysis, the most common statistical test used in comparing
1385 models in cross-validation is the paired t-test on the fold-level accuracy differences, referred
1386 to as the “resampled paired t-test” (Nadeau & Bengio, 2003). Given a vector \mathbf{D} of J fold-level
1387 performance differences between two models, we test the null hypothesis $\mu = 0$, where μ is
1388 the true expected performance difference between the two models. The paired t-test utilizes
1389 the following statistic:

1390

1391

$$T = \frac{\bar{D}}{\sqrt{\frac{1}{J} S^2}}$$

1392 where \bar{D} is the mean of vector \mathbf{D} , $S^2 = \frac{1}{J-1} \sum_j (D_j - \bar{D})^2$ is the sample variance of vector \mathbf{D} ,

1393

1394 The paired t-test assumes independence among the elements of vector \mathbf{D} , an assumption that
1395 is violated under cross-validation. In K -fold cross-validation, test folds are mutually disjoint
1396 across iterations, but the training sets overlap substantially (for $K > 2$), and the test fold in
1397 one iteration contributes to the training set in another iteration. In repeated K -fold cross-
1398 validation, the folds are redefined, but the cross-dependence remains. Similarly, in Monte
1399 Carlo cross-validation, where the dataset is repeatedly split into training and test sets, there
1400 will be overlap of both training and test sets across random dataset splits.

1401

1402 Because the correlations among fold-level statistics are positive, the paired t-test is expected
1403 to have elevated false positive rates, as demonstrated empirically in Fig. 5. A more detailed
1404 discussion of this issue can be found in Supplementary Methods S3.

1405

1406 4.5.2 Wilcoxon signed-rank test

1407 As demonstrated in our meta-analysis, the second most common statistical test used in
1408 comparing models in cross-validation is the Wilcoxon signed-rank or Wilcoxon rank-sum
1409 test. The Wilcoxon rank-sum test is also referred to as the Mann–Whitney U test, and is the
1410 unpaired version of the Wilcoxon signed-rank test. The Wilcoxon signed-rank test is a non-
1411 parametric alternative to the paired t-test, designed to test whether the median of paired
1412 differences is zero (Wilcoxon, 1945; Sidney, 1957).

1413

1414 Similar to the paired t-test, the Wilcoxon signed-rank test takes in a vector of paired
1415 differences and ranks the absolute differences between the J paired observations. The test
1416 then computes a test statistic based on the sum of these ranks, weighted by the sign of the

1417 original differences. The resulting test statistic is then compared against a reference
1418 distribution to obtain a p-value (Japkowicz & Shah, 2011).

1419
1420 Like the paired t-test, the Wilcoxon signed-rank test assumes independence among the paired
1421 differences. Because the fold-level statistics are positively correlated, the Wilcoxon signed-
1422 rank test is expected to have elevated false positive rates, as demonstrated empirically in Fig.
1423 5.

1424 1425 4.5.3 Sign-flip permutation test

1426 We found that there are two different permutation tests used in literature, but both are invalid.
1427 One permutation test is invalid not because of fold dependence but because the wrong
1428 variables are permuted. We excluded this permutation test from our meta-analysis. For
1429 details, see Supplementary Results.

1430
1431 We will describe the second permutation test here. Once again, the null hypothesis is $\mu = 0$,
1432 where μ is the true expected difference between the two models. Given a vector \mathbf{D} of J fold-
1433 level performance differences between two models, let \bar{D} be the mean of vector \mathbf{D} . A
1434 permutation of the data generates a null distribution: For each iteration, the sign of each
1435 element of vector \mathbf{D} is flipped independently with probability 0.5, thus creating a new vector
1436 \mathbf{D}_0 . The mean of \mathbf{D}_0 is then computed, contributing a single null value to the null
1437 distribution. After many iterations (e.g., 10,000), the original statistic \bar{D} is compared against
1438 the null distribution to generate a p-value.

1439
1440 This permutation test wrongly assumes the elements of vector \mathbf{D} are independent. Therefore,
1441 similar to the naïve t-test, we expect the permutation test to yield higher FPR, as
1442 demonstrated empirically in Fig. S5. A more detailed discussion of this issue can be found in
1443 Supplementary Methods S8.

1444 1445 4.5.4 DeLong's test

1446 The DeLong's test is a statistical procedure specific to model comparison based on a
1447 Receiver Operating Characteristic (ROC) curve (DeLong et al., 1988). The test only applies
1448 to binary classification tasks, in which the model predicts one of two possible outcomes (e.g.,
1449 spam vs. non-spam, disease vs. no disease). In such tasks, classifier performance can be
1450 summarized by an ROC curve, which plots the True Positive Rate (TPR) against the False
1451 Positive Rate (FPR) across different decision thresholds. To compare two ROC curves, the
1452 Area Under the ROC Curve (AUC) is used as a scalar performance metric, with larger AUC
1453 values indicating better performance.

1454
1455 DeLong's test evaluates whether the difference in AUCs between two algorithms is
1456 statistically significant. Specifically, the test statistic is constructed as the difference between
1457 the two estimated AUCs divided by an estimate of its standard error, yielding an asymptotic
1458 z-statistic. The AUC estimator is a function of the pairwise comparisons between positive
1459 and negative samples and can be expressed as a normalized Mann–Whitney U-statistic.
1460 DeLong's test estimates the 2×2 covariance matrix of the two AUC estimators using an
1461 influence-function-based approach for vectors of U-statistics (DeLong et al., 1988). The joint
1462 asymptotic normality of the AUC estimators justifies the use of a normal approximation to
1463 compute p-values.

1464
1465 Although the variance estimation in DeLong's test is mathematically more involved than that
1466 of simpler tests, the U-statistic variance derivation underlying DeLong's test assumes that

1467 predictions are independent across data samples, which is violated under cross-validation. We
1468 also note that the DeLong's test is based on sample-level statistics, as opposed to fold-
1469 averaged statistics (see Methods Section 4.1). As such, we expect the DeLong's test to have
1470 elevated FPR, as empirically shown in Fig. S1.

1471

1472 4.5.5 Corrected resampled paired t-test

1473 The corrected resampled paired t-test was proposed to adjust the variance estimator to reflect
1474 the positive correlation between cross-validation folds (Nadeau & Bengio, 2003). Although
1475 originally developed for Monte Carlo cross-validation, the corrected paired t-test has also
1476 been applied to repeated K-fold cross-validation (Bouckaert & Frank, 2004).

1477

1478 Similar to previous tests, the corrected resampled t-test operates on the vector \mathbf{D} of J fold-
1479 level performance differences between two models. Again, we wanted to test the null
1480 hypothesis $\mu = 0$, where μ is the true expected difference between the two models. The
1481 corrected resampled t-test assumes that the correlation between folds is equal to the fraction
1482 of the full dataset used in the test fold, resulting in the following statistic:

1483

$$1485 \quad T = \frac{\bar{D}}{\sqrt{\left(\frac{1}{J} + \frac{N_2}{N_1}\right) S^2}},$$

1484

1486 where \bar{D} is the mean of vector \mathbf{D} , $S^2 = \frac{1}{J-1} \sum_j (D_j - \bar{D})^2$ is the sample variance of vector \mathbf{D} ,
1487 and D_j refers to the j -th element of vector \mathbf{D} ; N_1 and N_2 are the number of training and test
1488 samples, respectively, in a particular split of the dataset into training and test sets. The
1489 statistic T is assumed to follow a Student's t distribution with $J-1$ degrees of freedom, from
1490 which a p-value can be computed.

1491

1492 A more detailed discussion of the corrected resampled t-test can be found in Supplementary
1493 Methods S3. One problem with the corrected resampled t-test is the assumption that the
1494 correlation between cross-validation folds is solely due to the overlap in the test set. In
1495 reality, the correlation likely depends on the complex interaction between the machine
1496 learning algorithms and the dataset being analyzed. We also note that some studies applied
1497 the corrected resampled t-test wrongly, resulting in a high FPR (Fig. S5; Supplementary
1498 Methods S9).

1499

1500 4.5.6 5×2 paired t-test and 5×2 paired F-test

1501 The 5×2 t-test (Dietterich, 1998) and 5×2 F-test (Alpaydin, 1999) involve 5 repetitions of 2-
1502 fold cross-validation, hence the name 5×2 . For each repetition $r \in \{1, 2, 3, 4, 5\}$, the data is
1503 randomly split into two halves A and B. Both models are trained on set A and evaluated on
1504 set B, resulting in performance difference D_{Br} ; likewise, both models are trained on set B and
1505 evaluated on set A, resulting in performance difference D_{Ar} .

1506

1507 The 5×2 t-test utilizes the following statistic (Dietterich, 1998):

1508

$$1511 \quad T = \frac{D_{A1}}{\sqrt{\frac{1}{5} \sum_{r=1}^5 S_r^2}}$$

1509 where $S_r^2 = (D_{Ar} - D_{Br})^2 / 2$. The statistic is assumed to follow a Student's t distribution
1510 with 5 degrees of freedom, from which a p-value can be computed. Supplementary Methods

1512 S4.1 explains why the 5×2 t-test will have a much lower false positive rate than the naïve t-
1513 test. However, the numerator in the T statistic D_{A1} only uses a single accuracy difference
1514 from the first fold of the first 2-fold cross-validation, thus discarding information from the
1515 other folds.

1516
1517 To address the inefficiency of the 5×2 t-test, Alpaydin (Alpaydin, 1999) proposed the 5×2
1518 paired F-test, which uses the following statistic:
1519

$$1521 \quad F = \frac{\sum_{r=1}^5 (D_{Ar}^2 + D_{Br}^2)}{2 \sum_{r=1}^5 S_r^2}$$

1520
1522 where $S_r^2 = (D_{Ar} - D_{Br})^2 / 2$ (same as the 5×2 t-test). The statistic is assumed to follow the
1523 F-distribution with 10 and 5 degrees of freedom in the numerator and denominator
1524 respectively, from which a p-value can be defined. Supplementary Methods S4.2 explains
1525 why the 5×2 F-test will have a much lower false positive rate than the naïve t-test. Intuitively,
1526 since the numerator uses fold-level differences from all 5 repetitions of 2-fold cross-
1527 validation (as opposed to just a single fold in the 5×2 t-test), the 5×2 F-test might be a more
1528 sensitive test than the 5×2 t-test (Alpaydin, 1999).

1529
1530 Indeed, our results suggest that both the 5×2 t-test and 5×2 F-test reliably control FPR, and
1531 that the 5×2 F-test exhibits higher power than the 5×2 t-test. However, both tests had lower
1532 power than SHARP.

1533 1534 4.5.7 Empirical test of differences

1535 Given a vector \mathbf{D} of J fold-level performance differences between two models, suppose the
1536 number of positive entries is greater than the number of negative entries. The empirical test
1537 of differences then computes the p-value as the fraction of negative entries multiplied by two.

1538
1539 On the other hand, suppose the number of positive entries is less than the number of negative
1540 entries. The empirical test of differences then computes the p-value as the fraction of positive
1541 entries multiplied by two.

1542
1543 The empirical test of differences has been utilized by a few brain imaging studies (Dhamala
1544 et al., 2021; Parkes et al., 2021a; Parkes et al., 2021b), which referred to it as the exact test of
1545 differences. Our results suggest that the empirical test of differences reliably controls false
1546 positives (Fig. 6) while having the lowest power among statistical tests that account for fold
1547 dependence (Fig. 7). Supplementary Methods S5 provides a theoretical explanation for these
1548 findings.

1549 1550 4.5.8 Three bootstrap variants

1551 Bootstrap is a common technique for estimating confidence intervals (Efron & Tibshirani,
1552 1994; DiCiccio & Efron, 1996). There are different bootstrapping variants for cross-
1553 validation (Raschka, 2018; Cai et al., 2025). Our current meta-analysis excluded all studies
1554 that only reported a confidence interval (PRISMA criterion 4; Methods Section 4.2). As such,
1555 all studies using bootstrap were also excluded. However, we note that only one study
1556 described their bootstrapping procedure in sufficient details for us to evaluate its validity
1557 (Peneder et al., 2021).

1558
1559 Below we outlined three broad bootstrap strategies that we evaluated in the current study.
1560 The three strategies were not meant to be exhaustive but were aimed to cover the good and

1561 bad variants of bootstrap. First, we implemented the bootstrap procedure from Raschka
1562 (2018), since it coincided with the exemplary study in our meta-analysis that clearly
1563 described the procedure (Peneder et al., 2021). More specifically, given a dataset of N
1564 samples, we bootstrapped (sampled with replacement) N samples from the dataset, which
1565 served as the training set. We note that on average, 63% of the training samples will be
1566 unique. The test set comprised samples that were not included in the training set.

1567
1568 To compare two models, the bootstrapped training set was used to train each model and the
1569 trained models were evaluated on the test set, yielding a single bootstrapped difference value.
1570 The entire bootstrapping procedure was repeated M times, and the 95% confidence interval
1571 was constructed from the M performance difference estimates based on the 2.5 and 97.5
1572 percentiles. If the 95% confidence interval did not cover zero, we considered the difference to
1573 be statistically significant. We referred to this bootstrapping procedure as “bootstrap-orig”.
1574 For the evaluation of FPR and power, we followed the procedure in Section 4.4 with M being
1575 set to 300.

1576
1577 The other two bootstrapping procedures performed bootstraps after normal cross-validation.
1578 For both approaches, given a vector \mathbf{D} of J fold-level performance differences, bootstrapping
1579 (sampling with replacement) was performed on the entries of vector \mathbf{D} , yielding a new vector
1580 \mathbf{D}^* of 1000 bootstrapped fold-level differences. An empirical test of differences (Methods
1581 Section 4.5.7) was applied to \mathbf{D}^* , which we referred to as bootstrap-ET. Alternatively, the
1582 naïve paired t-test (Section 4.5.1) was applied to \mathbf{D}^* , which we referred to as bootstrap-t.

1583
1584 Our results suggest that bootstrap-orig and bootstrap-ET reliably control FPR, but bootstrap-t
1585 has high FPR (Fig. S5). However, both bootstrap-orig and bootstrap-ET had markedly worse
1586 power than SHARP, with bootstrap-ET being the least powerful (Fig. S7).

1587 1588 **4.6 Split-Half Repeated (SHARP) test**

1589 The challenge of inference on performance differences comes down to estimating the
1590 standard error of the sample mean difference while accounting for the dependence between
1591 the J differences. Perhaps the clearest illustration of the problem is to consider a single
1592 application of K -fold cross-validation (discussed in Section 2.6): The J entries in the vector \mathbf{D}
1593 of fold-level differences are correlated, and for Gaussian data, there are two sufficient
1594 statistics (sample mean and sample variance), but three unknown parameters (mean, variance,
1595 and correlation). Thus, the model is overparameterized and, specifically, the variance and
1596 correlation cannot both be estimated without restrictive assumptions. The same fundamental
1597 statistical ambiguity applies to Monte Carlo cross-validation.

1598
1599 In their seminal study, Bengio and Grandvalet (Bengio & Grandvalet, 2004) considered a
1600 single instance of K -fold cross-validation and worked from sample-level prediction errors
1601 instead of fold-averaged performance (see Methods Section 4.1). They showed that there is
1602 no universal, distribution-independent unbiased estimator of the standard error. In particular,
1603 for Gaussian data, the maximum likelihood estimator does not exist. The covariance structure
1604 in their setting (based on sample-level prediction performance) is equivalent to the
1605 covariance structure of fold-averaged prediction performance under repeated K -fold cross-
1606 validation, so their findings generalize to repeated K -fold cross-validation.

1607
1608 The key idea behind the SHARP test is to create a situation where independent data are
1609 generated that will allow estimation of the variance parameters σ^2 and ρ . To achieve this, we
1610 randomly divide the dataset into two disjoint halves A and B. For each machine learning

1611 model, we then perform K-fold cross-validation in subsets A and B separately. We then
1612 average the results across the K-fold cross-validation, resulting in a pair of model
1613 performance differences, D_{A1} and D_{B1} respectively. Alternatively, one iteration of Monte
1614 Carlo cross-validation can be done within each subset A and B, likewise producing a pair of
1615 performance D_{A1} and D_{B1} based on the held-out testing data in each half.

1616
1617 We note that D_{A1} and D_{B1} are independent since subsets A and B are non-overlapping. For
1618 the analyses in the current study, this process is repeated J times, resulting in two vectors
1619 $\mathbf{D}_A = [D_{A1}, \dots, D_{AJ}]$ and $\mathbf{D}_B = [D_{B1}, \dots, D_{BJ}]$. We note that while D_{Aj} and D_{Bj} (from the j -th
1620 iteration) are independent, the entries within \mathbf{D}_A (and \mathbf{D}_B), and D_{Aj} and $D_{Bj'}$ for $j \neq j'$ are
1621 dependent with a common correlation. Therefore, we have J pairs of independent
1622 observations and three unknowns: mean performance μ , the variance σ^2 of each entry in \mathbf{D}_A
1623 and \mathbf{D}_B , and correlation ρ among all non-paired entries in \mathbf{D}_A and \mathbf{D}_B . Given \mathbf{D}_A and \mathbf{D}_B , we
1624 can estimate all three unknowns to perform the statistical test.

1625
1626 We define \bar{D} to be the grand average of the two vectors \mathbf{D}_A and \mathbf{D}_B . We seek to test the null
1627 hypothesis of equal expected performance:

$$H_0: \mu = 0, \text{ where } \mu = E[\bar{D}]$$

1630
1631 We show that the SHARP estimator of performance difference

$$\bar{D} = (\bar{D}_A + \bar{D}_B)/2,$$

1634
1635 where \bar{D}_A and \bar{D}_B are the respective sample means, is the generalized least-squares estimator
1636 of μ (Supplementary Methods S7.2), with variance

$$\text{Var}(\bar{D}) = \sigma^2 \left(\frac{1}{2J} + \frac{J-1}{J} \rho \right)$$

1639
1640 and so inference reduces to estimating σ^2 and ρ . To estimate σ^2 and ρ , we considered method-
1641 of-moments (Supplementary Methods S7.3), maximum likelihood (Supplementary Methods
1642 S7.4) or restricted maximum likelihood (Supplementary Methods S7.5), followed by a Wald
1643 test. We also considered likelihood-based tests including score test (Supplementary Methods
1644 S7.6) and likelihood ratio test (Supplementary Methods S7.7).

1645
1646 Based on simulated data (Supplementary Methods S7.8), we found that the score test
1647 provided the best control of FPR (Fig. S11). Therefore, all results in the current study utilized
1648 the score test. For the score test, the Z-score is a ratio of \bar{D} to a standard error computed using
1649 $\text{Var}(\bar{D})$ above, but with estimates $\hat{\sigma}_0^2$ and $\hat{\rho}_0$ obtained by optimizing the Gaussian likelihood
1650 assuming the null hypothesis is true, i.e., $\mu = 0$.

1651
1652 More details about the SHARP test can be found in Supplementary Methods S7. For the
1653 purpose of evaluating FPR and statistical power (Methods Section 4.4), we considered two
1654 settings: (1) $K = 5$, $J = 60$ and (2) $K = 10$, $J = 30$. Here, J was chosen to ensure convergence
1655 of the SHARP test.

1656

1657 **4.7 Ethics and data availability**

1658 Use of de-identified data from UKB datasets is approved by the National University of
1659 Singapore (NUS) Institutional Review Board (IRB).

1660

1661 This study used publicly available data from the UK Biobank

1662 (<https://www.ukbiobank.ac.uk/>), EMNIST

1663 (<https://www.kaggle.com/datasets/crawford/emnist>), Covertypes

1664 (<https://archive.ics.uci.edu/dataset/31/covertime>) and KEGG Metabolic Pathway

1665 (<https://archive.ics.uci.edu/dataset/220/kegg+metabolic+relation+network+directed>).

1666

1667 **4.8 Code availability**

1668 Code for this study can be found here (GITHUB_LINK). Co-authors (TZ and HL) reviewed
1669 each other's code before merging into the GitHub repository to reduce the chance of coding
1670 errors.

1671

1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683

Acknowledgements

Our research is supported by the NUS Yong Loo Lin School of Medicine (NUHSRO/2020/124/TMR/LOA), the Singapore National Medical Research Council (NMRC) LCG (OFLCG19May-0035), NMRC CTG-IIT (CTGIIT23jan-0001), NMRC OF-IRG (OFIRG24jan-0006; OFIRG24jul-0049), NMRC STaR (STaR20nov-0003), Singapore Ministry of Health (MOH) Centre Grant (CG21APR1009), the United States National Institutes of Health (R01MH133334 & 2R01MH120080) and the Singapore National Research Foundation (NRF) Investigatorship (NRFI10-2024-0014). This research has been conducted using the UK Biobank Resource under application number 25163. Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not reflect the views of the funders.

References

- 1684
1685 Alfaro-Almagro, F., Jenkinson, M., Bangerter, N. K., Andersson, J. L., Griffanti, L., Douaud,
1686 G., Sotiropoulos, S. N., Jbabdi, S., Hernandez-Fernandez, M., & Vallee, E. (2018).
1687 Image processing and Quality Control for the first 10,000 brain imaging datasets from
1688 UK Biobank. *Neuroimage*, *166*, 400–424.
- 1689 Ali, M. (2020). *PyCaret: An open source, low-code machine learning library in Python*.
1690 <https://www.pycaret.org>
- 1691 Alpaydin, E. (1999). Combined 5×2 cv F test for comparing supervised classification
1692 learning algorithms. *Neural Computation*, *11*(8), 1885–1892.
- 1693 André, P., Heitz, C., Christodoulou, E., Reinke, A., Sudre, C. H., Antonelli, M., Godau, P.,
1694 Cardoso, M. J., Gilson, A., Montcel, S. T. du, Varoquaux, G., Maier-Hein, L., &
1695 Colliot, O. (2026). *Performance uncertainty in medical image analysis: A large-scale*
1696 *investigation of confidence intervals* (arXiv:2601.17103). arXiv.
1697 <https://doi.org/10.48550/arXiv.2601.17103>
- 1698 Bates, S., Hastie, T., & Tibshirani, R. (2024). Cross-Validation: What Does It Estimate and
1699 How Well Does It Do It? *Journal of the American Statistical Association*, *119*(546),
1700 1434–1445. <https://doi.org/10.1080/01621459.2023.2197686>
- 1701 Bengio, Y., & Grandvalet, Y. (2004). No unbiased estimator of the variance of k-fold cross-
1702 validation. *Journal of Machine Learning Research*, *5*(Sep), 1089–1105.
- 1703 Bergstra, J., Bardenet, R., Bengio, Y., & Kégl, B. (2011). Algorithms for hyper-parameter
1704 optimization. *Advances in Neural Information Processing Systems*, *24*.
1705 [https://proceedings.neurips.cc/paper/4443-algorithms-for-hyper-parameter-](https://proceedings.neurips.cc/paper/4443-algorithms-for-hyper-parameter-optimization)
1706 [optimization](https://proceedings.neurips.cc/paper/4443-algorithms-for-hyper-parameter-optimization)
- 1707 Bergstra, J., Yamins, D., & Cox, D. (2013). Making a science of model search:
1708 Hyperparameter optimization in hundreds of dimensions for vision architectures.
1709 *International Conference on Machine Learning*, 115–123.
1710 <https://proceedings.mlr.press/v28/bergstra13.html>
- 1711 Blackard, J. (1998). *Coverttype* [Dataset]. UCI Machine Learning Repository.
1712 <https://doi.org/10.24432/C50K5N>
- 1713 Bouckaert, R. R., & Frank, E. (2004). Evaluating the Replicability of Significance Tests for
1714 Comparing Learning Algorithms. In H. Dai, R. Srikant, & C. Zhang (Eds.), *Advances*
1715 *in Knowledge Discovery and Data Mining* (Vol. 3056, pp. 3–12). Springer Berlin
1716 Heidelberg. https://doi.org/10.1007/978-3-540-24775-3_3
- 1717 Bzdok, D., Thieme, A., Levkovskyy, O., Wren, P., Ray, T., & Reddy, S. (2024). Data science
1718 opportunities of large language models for neuroscience and biomedicine. *Neuron*,
1719 *112*(5), 698–717.
- 1720 Cai, B., Luo, Y., Guo, X., Pellegrini, F., Pang, M., De Moor, C., Shen, C., Charu, V., & Tian,
1721 L. (2025). Bootstrapping the cross-validation estimate. *The Annals of Applied*
1722 *Statistics*, *19*(4), 2981–3002.
- 1723 Carrasco-Zanini, J., Pietzner, M., Davitte, J., Surendran, P., Croteau-Chonka, D. C., Robins,
1724 C., Torralbo, A., Tomlinson, C., Grünschläger, F., & Fitzpatrick, N. (2024).
1725 Proteomic signatures improve risk prediction for common and rare diseases. *Nature*
1726 *Medicine*, *30*(9), 2489–2498.
- 1727 Chopra, S., Dhamala, E., Lawhead, C., Ricard, J. A., Orchard, E. R., An, L., Chen, P., Wulan,
1728 N., Kumar, P., Rubenstein, A., Moses, J., Chen, L., Levi, P., Holmes, A., Aquino, K.,
1729 Fornito, A., Harpaz-Rotem, I., Germine, L. T., Baker, J. T., ... Holmes, A. J. (2024).
1730 Generalizable and replicable brain-based predictions of cognitive functioning across
1731 common psychiatric illness. *Science Advances*, *10*(45), eadn1862.
1732 <https://doi.org/10.1126/sciadv.adn1862>

- 1733 Christodoulou, E., Reinke, A., Andr e, P., Godau, P., Kalinowski, P., Houhou, R., Erkan, S.,
1734 Sudre, C. H., Burgos, N., Boutaj, S., Loizillon, S., Solal, M., Cheplygina, V., Heitz,
1735 C., Kozubek, M., Antonelli, M., Rieke, N., Gilson, A., Mayer, L. D., ... Maier-Hein,
1736 L. (2025). *False Promises in Medical Imaging AI? Assessing Validity of*
1737 *Outperformance Claims* (arXiv:2505.04720). arXiv.
1738 <https://doi.org/10.48550/arXiv.2505.04720>
- 1739 Christodoulou, E., Reinke, A., Houhou, R., Kalinowski, P., Erkan, S., Sudre, C. H., Burgos,
1740 N., Boutaj, S., Loizillon, S., Solal, M., Rieke, N., Cheplygina, V., Antonelli, M.,
1741 Mayer, L. D., Tizabi, M. D., Cardoso, M. J., Simpson, A., J ager, P. F., Kopp-
1742 Schneider, A., ... Maier-Hein, L. (2024). Confidence Intervals Uncovered: Are We
1743 Ready for Real-World Medical Imaging AI? In M. G. Linguraru, Q. Dou, A. Feragen,
1744 S. Giannarou, B. Glocker, K. Lekadir, & J. A. Schnabel (Eds.), *Medical Image*
1745 *Computing and Computer Assisted Intervention – MICCAI 2024* (Vol. 15010, pp.
1746 124–132). Springer Nature Switzerland. [https://doi.org/10.1007/978-3-031-72117-](https://doi.org/10.1007/978-3-031-72117-5_12)
1747 [5_12](https://doi.org/10.1007/978-3-031-72117-5_12)
- 1748 Cohen, G., Afshar, S., Tapson, J., & Van Schaik, A. (2017). EMNIST: Extending MNIST to
1749 handwritten letters. *2017 International Joint Conference on Neural Networks*
1750 *(IJCNN)*, 2921–2926. <https://ieeexplore.ieee.org/abstract/document/7966217/>
- 1751 Collobert, R., Bengio, S., & Bengio, Y. (2001). A parallel mixture of SVMs for very large
1752 scale problems. *Advances in Neural Information Processing Systems, 14*.
1753 [https://proceedings.neurips.cc/paper_files/paper/2001/hash/36ac8e558ac7690b6f44e2](https://proceedings.neurips.cc/paper_files/paper/2001/hash/36ac8e558ac7690b6f44e2cb5ef93322-Abstract.html)
1754 [cb5ef93322-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2001/hash/36ac8e558ac7690b6f44e2cb5ef93322-Abstract.html)
- 1755 DeLong, E. R., DeLong, D. M., & Clarke-Pearson, D. L. (1988). Comparing the areas under
1756 two or more correlated receiver operating characteristic curves: A nonparametric
1757 approach. *Biometrics*, 837–845.
- 1758 Dem sar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of*
1759 *Machine Learning Research*, 7(Jan), 1–30.
- 1760 Dhamala, E., Jamison, K. W., Jaywant, A., Dennis, S., & Kuceyeski, A. (2021). Distinct
1761 functional and structural connections predict crystallised and fluid cognition in
1762 healthy adults. *Human Brain Mapping*, 42(10), 3102–3118.
1763 <https://doi.org/10.1002/hbm.25420>
- 1764 DiCiccio, T. J., & Efron, B. (1996). Bootstrap confidence intervals. *Statistical Science*, 11(3),
1765 189–228.
- 1766 Dietterich, T. G. (1998). Approximate statistical tests for comparing supervised classification
1767 learning algorithms. *Neural Computation*, 10(7), 1895–1923.
- 1768 Edgington, E., & Onghena, P. (2007). *Randomization tests*. Chapman and Hall/CRC.
1769 [https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&ide](https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&identifierValue=10.1201/9781420011814&type=googlepdf)
1770 [ntifierValue=10.1201/9781420011814&type=googlepdf](https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&identifierValue=10.1201/9781420011814&type=googlepdf)
- 1771 Efron, B., & Tibshirani, R. J. (1994). *An introduction to the bootstrap*. Chapman and
1772 Hall/CRC.
1773 [https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&ide](https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&identifierValue=10.1201/9780429246593&type=googlepdf)
1774 [ntifierValue=10.1201/9780429246593&type=googlepdf](https://api.taylorfrancis.com/content/books/mono/download?identifierName=doi&identifierValue=10.1201/9780429246593&type=googlepdf)
- 1775 Errington, T. M., Denis, A., Perfito, N., Iorns, E., & Nosek, B. A. (2021). Challenges for
1776 assessing replicability in preclinical cancer biology. *Elife*, 10, e67995.
- 1777 Fischl, B. (2012). FreeSurfer. *Neuroimage*, 62(2), 774–781.
- 1778 Greene, A. S., Gao, S., Scheinost, D., & Constable, R. T. (2018). Task-induced brain state
1779 manipulation improves prediction of individual traits. *Nature Communications*, 9(1),
1780 2807.

- 1781 Grinsztajn, L., Oyallon, E., & Varoquaux, G. (2022). Why do tree-based models still
1782 outperform deep learning on typical tabular data? *Advances in Neural Information*
1783 *Processing Systems*, 35, 507–520.
- 1784 He, T., An, L., Chen, P., Chen, J., Feng, J., Bzdok, D., Holmes, A. J., Eickhoff, S. B., & Yeo,
1785 B. T. (2022). Meta-matching as a simple framework to translate phenotypic predictive
1786 models from big to small data. *Nature Neuroscience*, 25(6), 795–804.
- 1787 Jafrasteh, B., Adeli, E., Pohl, K. M., Kuceyeski, A., Sabuncu, M. R., & Zhao, Q. (2025).
1788 Statistical variability in comparing accuracy of neuroimaging based classification
1789 models via cross validation. *Scientific Reports*, 15(1), 28745.
- 1790 Japkowicz, N., & Shah, M. (2011). *Evaluating learning algorithms: A classification*
1791 *perspective*. Cambridge University Press.
1792 [https://books.google.com/books?hl=en&lr=&id=VoWIIOKVzR4C&oi=fnd&pg=PR7](https://books.google.com/books?hl=en&lr=&id=VoWIIOKVzR4C&oi=fnd&pg=PR7&dq=Evaluating+Learning+Algorithms+(Cambridge)&ots=5z86ZNHZOK&sig=X4tJrh_j5TGme-EGiNqD_cf7O38)
1793 [&dq=Evaluating+Learning+Algorithms+\(Cambridge\)&ots=5z86ZNHZOK&sig=X4tJ](https://books.google.com/books?hl=en&lr=&id=VoWIIOKVzR4C&oi=fnd&pg=PR7&dq=Evaluating+Learning+Algorithms+(Cambridge)&ots=5z86ZNHZOK&sig=X4tJrh_j5TGme-EGiNqD_cf7O38)
1794 [rh_j5TGme-EGiNqD_cf7O38](https://books.google.com/books?hl=en&lr=&id=VoWIIOKVzR4C&oi=fnd&pg=PR7&dq=Evaluating+Learning+Algorithms+(Cambridge)&ots=5z86ZNHZOK&sig=X4tJrh_j5TGme-EGiNqD_cf7O38)
- 1795 Kapoor, S., & Narayanan, A. (2023). Leakage and the reproducibility crisis in machine-
1796 learning-based science. *Patterns*, 4(9). [https://www.cell.com/patterns/fulltext/S2666-](https://www.cell.com/patterns/fulltext/S2666-3899(23)00159-9)
1797 [3899\(23\)00159-9](https://www.cell.com/patterns/fulltext/S2666-3899(23)00159-9)
- 1798 Komer, B., Bergstra, J., & Eliasmith, C. (2014). Hyperopt-Sklearn: Automatic
1799 Hyperparameter Configuration for Scikit-Learn. *Scipy*, 32–37.
1800 [https://pub.curvenote.com/0190828e-e1d1-7a2c-8d33-0d909287203f/public/komer-](https://pub.curvenote.com/0190828e-e1d1-7a2c-8d33-0d909287203f/public/komer-20a7edaa18d743190b2a53628c177e8f.pdf)
1801 [20a7edaa18d743190b2a53628c177e8f.pdf](https://pub.curvenote.com/0190828e-e1d1-7a2c-8d33-0d909287203f/public/komer-20a7edaa18d743190b2a53628c177e8f.pdf)
- 1802 Mann, H. B., & Whitney, D. R. (1947). On a test of whether one of two random variables is
1803 stochastically larger than the other. *The Annals of Mathematical Statistics*, 50–60.
- 1804 Mansour L, S., Tian, Y., Yeo, B. T. T., Cropley, V., & Zalesky, A. (2021). High-resolution
1805 connectomic fingerprints: Mapping neural identity and behavior. *NeuroImage*, 229,
1806 117695. <https://doi.org/10.1016/j.neuroimage.2020.117695>
- 1807 Moor, M., Banerjee, O., Abad, Z. S. H., Krumholz, H. M., Leskovec, J., Topol, E. J., &
1808 Rajpurkar, P. (2023). Foundation models for generalist medical artificial intelligence.
1809 *Nature*, 616(7956), 259–265. <https://doi.org/10.1038/s41586-023-05881-4>
- 1810 Muhammad Naeem, S. A. (2011). *KEGG Metabolic Relation Network (Directed)* [Dataset].
1811 UCI Machine Learning Repository. <https://doi.org/10.24432/C5CK52>
- 1812 Nadeau, C., & Bengio, Y. (2003). Inference for the Generalization Error. *Machine Learning*,
1813 52(3), 239–281. <https://doi.org/10.1023/A:1024068626366>
- 1814 Newcombe, R. G. (1998). Two-sided confidence intervals for the single proportion:
1815 Comparison of seven methods. *Statistics in Medicine*, 17(8), 857–872.
1816 [https://doi.org/10.1002/\(SICI\)1097-0258\(19980430\)17:8<253C857::AID-](https://doi.org/10.1002/(SICI)1097-0258(19980430)17:8<253C857::AID-SIM777>2.3.CO;2-E)
1817 [SIM777>253E3.0.CO;2-E](https://doi.org/10.1002/(SICI)1097-0258(19980430)17:8<253C857::AID-SIM777>2.3.CO;2-E)
- 1818 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.
1819 *Science*, 349(6251), aac4716. <https://doi.org/10.1126/science.aac4716>
- 1820 Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D.,
1821 Shamseer, L., Tetzlaff, J. M., Akl, E. A., & Brennan, S. E. (2021). The PRISMA 2020
1822 statement: An updated guideline for reporting systematic reviews. *Bmj*, 372.
1823 <https://www.bmj.com/content/372/bmj.n71.short>
- 1824 Parkes, L., Moore, T. M., Calkins, M. E., Cieslak, M., Roalf, D. R., Wolf, D. H., Gur, R. C.,
1825 Gur, R. E., Satterthwaite, T. D., & Bassett, D. S. (2021a). Network controllability in
1826 transmodal cortex predicts positive psychosis spectrum symptoms. *Biological*
1827 *Psychiatry*, 90(6), 409–418.
- 1828 Parkes, L., Moore, T. M., Calkins, M. E., Cook, P. A., Cieslak, M., Roalf, D. R., Wolf, D. H.,
1829 Gur, R. C., Gur, R. E., & Satterthwaite, T. D. (2021b). Transdiagnostic dimensions of

- 1830 psychopathology explain individuals' unique deviations from normative
1831 neurodevelopment in brain structure. *Translational Psychiatry*, *11*(1), 232.
- 1832 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
1833 Prettenhofer, P., Weiss, R., & Dubourg, V. (2011). Scikit-learn: Machine learning in
1834 Python. *The Journal of Machine Learning Research*, *12*, 2825–2830.
- 1835 Peneder, P., Stütz, A. M., Surdez, D., Krumbholz, M., Semper, S., Chicard, M., Sheffield, N.
1836 C., Pierron, G., Lapouble, E., & Tötzl, M. (2021). Multimodal analysis of cell-free
1837 DNA whole-genome sequencing for pediatric cancers with low mutational burden.
1838 *Nature Communications*, *12*(1), 3230.
- 1839 Perez-Lopez, R., Ghaffari Laleh, N., Mahmood, F., & Kather, J. N. (2024). A guide to
1840 artificial intelligence for cancer researchers. *Nature Reviews Cancer*, *24*(6), 427–441.
1841 <https://doi.org/10.1038/s41568-024-00694-7>
- 1842 Rajpurkar, P., Chen, E., Banerjee, O., & Topol, E. J. (2022). AI in health and medicine.
1843 *Nature Medicine*, *28*(1), 31–38.
- 1844 Raschka, S. (2018). *Model Evaluation, Model Selection, and Algorithm Selection in Machine*
1845 *Learning* (arXiv:1811.12808). arXiv. <https://doi.org/10.48550/arXiv.1811.12808>
- 1846 Shannon, P., Markiel, A., Ozier, O., Baliga, N. S., Wang, J. T., Ramage, D., Amin, N.,
1847 Schwikowski, B., & Ideker, T. (2003). Cytoscape: A software environment for
1848 integrated models of biomolecular interaction networks. *Genome Research*, *13*(11),
1849 2498–2504.
- 1850 Sidney, S. (1957). Nonparametric statistics for the behavioral sciences. *The Journal of*
1851 *Nervous and Mental Disease*, *125*(3), 497.
- 1852 Varoquaux, G., Raamana, P. R., Engemann, D. A., Hoyos-Idrobo, A., Schwartz, Y., &
1853 Thirion, B. (2017). Assessing and tuning brain decoders: Cross-validation, caveats,
1854 and guidelines. *NeuroImage*, *145*, 166–179.
- 1855 Wagner, S. J., Reisenbüchler, D., West, N. P., Niehues, J. M., Zhu, J., Foersch, S.,
1856 Veldhuizen, G. P., Quirke, P., Grabsch, H. I., & van den Brandt, P. A. (2023).
1857 Transformer-based biomarker prediction from colorectal cancer histology: A large-
1858 scale multicentric study. *Cancer Cell*, *41*(9), 1650–1661.
- 1859 Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics Bulletin*, *1*(6),
1860 80–83.
- 1861 Wright, S. (1921). Correlation and causation. *Journal of Agricultural Research*, *20*(7), 557.
- 1862 Wulan, N., An, L., Zhang, C., Kong, R., Chen, P., Bzdok, D., Eickhoff, S. B., Holmes, A. J.,
1863 & Yeo, B. T. (2024). Translating phenotypic prediction models from big to small
1864 anatomical MRI data using meta-matching. *Imaging Neuroscience*, *2*, 1–21.
- 1865 Yeo, B. T., Sabuncu, M. R., Vercauteren, T., Holt, D. J., Amunts, K., Zilles, K., Golland, P.,
1866 & Fischl, B. (2010). Learning task-optimal registration cost functions for localizing
1867 cytoarchitecture and function in the cerebral cortex. *IEEE Transactions on Medical*
1868 *Imaging*, *29*(7), 1424–1441.
- 1869 Yoo, S.-K., Fitzgerald, C. W., Cho, B. A., Fitzgerald, B. G., Han, C., Koh, E. S., Pandey, A.,
1870 Sfreddo, H., Crowley, F., & Korostin, M. R. (2025). Prediction of checkpoint
1871 inhibitor immunotherapy efficacy for cancer using routine blood tests and clinical
1872 data. *Nature Medicine*, *31*(3), 869–880.
- 1873
1874