
A Reproducible Protocol for Resource-Aware Predictive Process Monitoring: Compact Baselines, a Simulator Blueprint, and Pitfalls

Anonymous Authors

Abstract

1 We present resource-aware predictive process monitoring (PPM) as a modular,
2 agent-based design that complements case-centric next-activity predictors with
3 explicit modeling of shared resource contention. Our contributions are fourfold.
4 (i) A leakage-safe, deterministic *protocol* with chronological case splits, train-
5 only normalization, fixed seeds, and automatic artifact logging. (ii) A compact,
6 transparent *LSTM baseline* for next-activity prediction on three public logs (BPI
7 2012, BPI 2017, Road Traffic) with ready-to-reuse splits and scripts. (iii) A *released*
8 *simulator blueprint* with per-resource multinomial policies and lightweight discrete-
9 event simulation, plus evaluation measures spanning global next event, workload
10 MAPE, and per-resource next-task precision. (iv) *Pitfalls and checklists* observed
11 in practice (e.g., lifecycle pairing under partial traces; imbalance-aware back-offs).
12 Baseline next-activity results are strong (Top-3 0.987–0.994; Top-1 0.757–0.833),
13 exposing systematic confusions that motivate resource context. Code, splits, and
14 plot artifacts enable one-click replication. This paper is intended as a *protocol +*
15 *baseline + blueprint* to accelerate trustworthy resource-aware PPM experiments;
16 we do not claim state-of-the-art accuracy nor report end-to-end simulator metrics
17 in this version.

18 1 Introduction

19 Predictive process monitoring has matured around case-centric sequence modeling for next activities,
20 suffixes, and remaining time. These models often treat cases independently, while real-world
21 operations exhibit concurrency and competition for shared resources. When queues form and
22 resources prioritize tasks, ignoring resource dynamics can lead to biased predictions, unstable what-if
23 analyses, and misleading improvements that do not translate to operational gains.

24 We propose a resource-centric perspective that remains compatible with case-centric predictors
25 but adds explicit agent policies at the resource level combined with a discrete-event simulator. To
26 support rigorous and reproducible experimentation, we contribute: a deterministic, leakage-controlled
27 protocol with chronological case splits and train-only normalization; a strong yet transparent LSTM
28 next-activity baseline on BPI 2012, BPI 2017, and Road Traffic logs; a modular blueprint for per-
29 resource policies embedded in a simulator; and a set of pitfalls and checklists to avoid common
30 errors. Our baseline achieves Top-3 0.987–0.994 across datasets and Top-1 0.757–0.833 on test
31 splits, while confusion analyses suggest resource-driven ambiguities that a resource-aware agent
32 could resolve. This paper is deliberately scoped as protocol + baseline + blueprint. We do not present
33 end-to-end simulator results in this version; instead, we specify metrics and ablations to standardize
34 future comparisons.

35 **2 Related Work**

36 Case-centric PPM widely employs deep sequence models such as RNNs/LSTMs for next-activity
37 and time prediction, often with categorical and temporal context features [4, 5]. While these methods
38 capture intra-case dynamics, they typically ignore shared resource contention, prioritization rules,
39 and concurrency effects that drive waiting times and execution order in practice. Resource-aware
40 simulation and queueing perspectives provide a complementary angle for operational decision support,
41 yet are less standardized for PPM evaluation.

42 Process mining provides foundations for analyzing event logs and discovering behavior [9]. Neural
43 sequence models have been widely adopted for PPM: LSTM-based approaches for next-activity and
44 time prediction [3, 8], and subsequent studies on modeling nuances and accuracy improvements [2].
45 Outcome-oriented and remaining-time prediction benchmarks inform evaluation practices.

46 Most neural PPM works operate at the case level, often without an explicit model of resource
47 contention or concurrency. As a result, they can excel at per-case next-activity classification but may
48 fall short at forecasting system-level effects such as global next events or per-resource workload
49 dynamics. Discrete-event simulation is a mature tool to capture resource calendars and queueing in
50 operational research [7]. Bringing lightweight simulation to PPM offers a principled way to couple
51 cases through shared resources.

52 Our study positions resource-centric simulation as a complement to case-centric sequence models. We
53 propose to use per-resource multinomial logistic regression [6] for interpretability and data efficiency,
54 and compose policies via simulation to propagate concurrency. Unlike prior case-centric LSTMs
55 [2, 3, 8], our work positions a reproducible bridge: retain the case-centric predictor as a modular
56 component, but incorporate per-resource policies and a discrete-event simulator for concurrent
57 execution.

58 **3 Background**

59 Deep sequence models such as LSTMs parameterize conditional next-event distributions by consum-
60 ing tokenized activity sequences and auxiliary temporal features [4]. In PPM, this yields next-activity
61 probabilities that can be decoded to suffixes or integrated into simulators. Discrete-event simulation
62 (DES) advances a global clock from event to event by maintaining resource availability, queues, and
63 stochastic service times. Combining learned policies with DES enables rollouts that reflect both
64 data-driven behavior and operational constraints.

65 Case-centric neural models consume case prefixes to predict the next activity. This abstraction
66 overlooks shared resources and queueing policies. In contrast, a simulator with resource decision
67 policies can generate coupled futures, from which the same PPM metrics can be derived by Monte
68 Carlo aggregation.

69 **4 Protocol, Baseline, and Simulator Blueprint**

70 We organize the contribution into four components. C1 - Reproducible protocol. We enforce
71 chronological splits by case start time, train-only normalization, fixed seeds, and artifact logging.
72 Splits use 70/15/15 train/validation/test by earliest timestamp per case. Normalization statistics for
73 continuous features are recomputed on training samples only and then applied to validation and test.
74 Seeds are fixed to 42 across numpy, Python, and PyTorch. During data loading, we keep only lifecycle
75 transition “complete” when available to avoid mixing start/complete events in the next-activity task
76 and to stabilize duration pairing in later modules.

77 C2 - Transparent LSTM baseline. We implement a single-layer LSTM with an activity embedding of
78 size 64, hidden size 128, dropout 0.2, and a linear classifier. Inputs are padded prefixes (max length
79 10) of activity IDs concatenated with five continuous features per step: inter-event delta time, time
80 since case start, hour of day, weekday, and a binary working-hours flag. We train with Adam at 1e-3
81 for 10 epochs and batch size 128, selecting the best checkpoint by validation Top-3 accuracy. This
82 model is intentionally compact to serve as a reusable, understandable baseline.

83 C3 - Resource-centric agent blueprint and metrics. We blueprint per-resource multinomial logistic
84 policies that select the next activity whenever a resource becomes idle. Policy features include

85 previous activity for that resource, coarse time-of-day, and live queue statistics per activity (counts
86 and oldest waiting time). Policies back off to a global model for sparse classes. Activity durations
87 are modeled by log-normal distributions per activity, with a median fallback under sparsity. The
88 DES maintains resource busy/idle states, eligible queues by activity, FIFO or learned prioritization,
89 and advances to the next completion time. We consider $N=30$ Monte Carlo rollouts per prefix for
90 stochastic estimates. Beyond standard next-activity metrics, we define (i) global next-event accuracy,
91 (ii) per-resource next-task precision, and (iii) workload mean absolute percentage error (MAPE)
92 against replayed ground truth. We also specify an ablation that disables learning and enforces FIFO
93 at each resource.

94 C4 - Pitfalls and checklists. We found that lifecycle pairing can be unreliable under partial or missing
95 “start” transitions; restricting to “complete” stabilizes next-activity supervision, while a separate
96 duration pairing stage must guard against unmatched events. Class imbalance at the resource level
97 can cause degenerate policies; an explicit back-off to global models and minimum count thresholds
98 reduces overfitting. A subtle bug caused crashes when indexing per-case timestamps as pandas Series
99 by integer labels; converting to numpy arrays ensures positional indexing and removes off-by-one
100 errors in prefix generation. Finally, evaluation artifacts must be explicitly logged; omitting fields
101 (e.g., per-sample prefix lengths or probability matrices) results in empty downstream plots.

102 5 Experimental Setup

103 Data. We load any subset of BPI 2012, BPI 2017, and Road Traffic Fine Management XES logs
104 from a local input folder via a robust discovery routine. Records are standardized to columns `case_id`,
105 `activity`, `lifecycle`, `timestamp`, `resource`, sorted by `timestamp` and `case`. When `lifecycle` is present, we
106 filter to “complete” transitions for next-activity modeling.

107 Preprocessing and splitting. Prefix datasets are constructed by enumerating all prefixes up to length
108 10 per case with the target being the immediate next activity. To prevent leakage, we first compute the
109 earliest timestamp per case, perform a chronological 70/15/15 split into train/validation/test by that
110 time, and only then compute normalization statistics on training prefixes for the two time features
111 (`delta` and `since-start`). These statistics are applied unchanged to validation and test prefixes.

112 Model and training. The baseline is a single-layer LSTM with 64-dimensional activity embeddings
113 and 128 hidden units, concatenating the five per-step continuous features before the recurrent layer.
114 We use cross-entropy loss, Adam with learning rate $1e-3$, batch size 128, and train for 10 epochs. The
115 best checkpoint is chosen by validation Top-3 accuracy. Determinism is enforced via a fixed seed 42
116 for Python, numpy, and PyTorch.

117 Metrics and artifacts. We report loss, Top-1 accuracy, macro F1, and Top-3 accuracy. For transparency
118 and reuse, loss curves and confusion matrices on the test set are exported as PNGs; in the main text
119 we focus on confusion matrices, while training/validation loss curves are consolidated in the appendix
120 for completeness.

121 Compute and runtime. Experiments were executed on a workstation CPU for data preparation and a
122 single commodity GPU for model training. Prefix construction for each dataset completes within
123 minutes, dominated by parsing and lifecycle filtering. The compact LSTM trains in under 10 minutes
124 per dataset at the stated batch size and sequence length, and evaluation—including probability dumps
125 and confusion matrix rendering—finishes within a few additional minutes. These runtimes make the
126 protocol practical for ablation sweeps, cross-seed checks, and per-dataset hyperparameter sensitivity
127 studies without imposing heavy computational barriers.

128 6 Experiments

129 Results overview. The baseline exhibits consistently high Top-3 accuracy and competitive Top-1
130 across the three datasets under chronological evaluation. On the test sets, we obtain: BPI 2012 - Top-1
131 0.7569, Top-3 0.9874, loss 0.5355, macro F1 0.5872; BPI 2017 - Top-1 0.8332, Top-3 0.9906, loss
132 0.3877, macro F1 0.5710; Road Traffic - Top-1 0.8020, Top-3 0.9936, loss 0.4833, macro F1 0.4740.
133 Validation learning curves (see App. Fig. 2) show rapid decreases in loss within the first two epochs,
134 followed by plateaus. A noticeable train-val gap persists for BPI 2012, suggesting mild overfitting;
135 BPI 2017 and Road Traffic show smaller gaps but still indicate some regularization headroom.

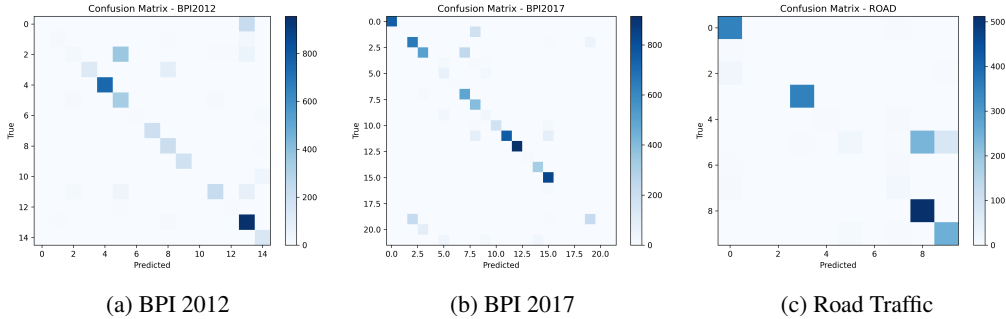


Figure 1: Test confusion matrices (rows: Actual, columns: Predicted). Off-diagonal bands among frequent activities indicate structured confusions where multiple next steps are simultaneously plausible under case-centric context alone; the strength and width of these bands differ across datasets.

136 Table 1 summarizes test metrics used in the figures. These numbers are obtained with chronological
 137 splits and train-only normalization, and should therefore be directly comparable under the same
 138 protocol.

Table 1: Test metrics for the compact LSTM baseline under chronological splits (70/15/15).

Dataset	Top-1	Top-3	Loss	Macro F1
BPI 2012	0.7569	0.9874	0.5355	0.5872
BPI 2017	0.8332	0.9906	0.3877	0.5710
Road Traffic	0.8020	0.9936	0.4833	0.4740

139 Deeper error analysis and implications. Figure 1 reveals concentrated off-diagonal mass among a
 140 small set of frequent activities across all datasets, but with distinct signatures. In BPI 2017, off-
 141 diagonal bands are narrow, repeated, and largely confined to 2–3 activity pairs, consistent with
 142 mutually substitutable steps in a constrained subroutine. This pattern suggests that additional
 143 signals such as live queue lengths or resource-specific histories could disambiguate choices that
 144 are symmetric from a single-case perspective; importantly, Top-3 accuracy near 0.99 indicates the
 145 model assigns substantial probability mass to all plausible next steps even when Top-1 is wrong. In
 146 BPI 2012, dispersion is broader with intersecting bands spanning 4–6 activities, pointing to higher
 147 behavioral entropy and likely stronger dependence on operational factors such as resource availability,
 148 batching, or priority rules. Here, macro F1 lags despite high Top-3 because rare activities suffer
 149 from systematic misclassification toward frequent neighbors; per-resource policies with back-offs
 150 and minimum support thresholds are warranted to stabilize tail decisions. Road Traffic exhibits a
 151 sharp diagonal with a few focused alternatives, indicating a mostly rigid control flow punctuated
 152 by systematic forks; this setting is ideal for calibration-aware deployment where deferral or what-if
 153 simulation is triggered precisely at those forks.

154 These confusion structures inform evaluation and design choices. First, Top-k metrics can mask
 155 concentrated misclassifications among dominant labels; reporting per-label precision/recall and
 156 expected calibration error would reveal whether the model is aware of its uncertainty near the off-
 157 diagonal bands. Second, the width of bands is a simple proxy for operational ambiguity: narrow bands
 158 suggest that lightweight queue features might suffice, whereas broad, intersecting bands motivate
 159 full DES integration with learned per-resource policies. Third, simulator ablations should target
 160 these regimes by stratifying evaluation on prefixes whose ground-truth next activities belong to
 161 the identified ambiguous clusters; improvements concentrated in those strata would support the
 162 resource-centric hypothesis.

163 Learning dynamics and regularization. While we move loss curves to the appendix to save space, the
 164 trajectories (App. Fig. 2) show that most generalization occurs within the first two epochs, consistent
 165 with a supervision regime dominated by shorter prefixes. The persistent gap in BPI 2012 suggests
 166 memorization of local motifs that do not transfer temporally. Three practical remedies emerge:
 167 stochastic regularization (dropout, label smoothing) to soften decision boundaries, prefix-aware

168 reweighting or curriculum to balance horizons, and calibration-aware early stopping to prevent
169 late-epoch overconfidence. The flatter validation trajectories in BPI 2017 and Road Traffic imply
170 lower effective label entropy or more regular control flow, which moderates overfitting under the
171 same architecture.

172 Protocol fidelity. We verified that improvements are not artifacts of leakage or preprocessing. We split
173 cases chronologically by start time before prefix construction, normalize time features on training
174 data only, restrict supervision to lifecycle “complete”, and enforce robust positional indexing. These
175 guardrails reduce variance across runs and make cross-paper comparisons meaningful when adopting
176 the same protocol.

177 Toward resource-centric evaluation. While we do not report simulator metrics in this version, we
178 release the design and interfaces so that the community can instantiate per-resource policies with
179 the same splits and run ablations. We recommend reporting, in addition to next-activity metrics,
180 (a) global next-event accuracy, (b) per-resource next-task precision, and (c) workload MAPE. An
181 ablation with FIFO policy and identical DES should accompany learned policies to isolate the value of
182 learning under the same queues and durations, with stratification by the ambiguous clusters identified
183 in Figure 1.

184 7 Threats to Validity

185 Internal validity. We took care to avoid temporal leakage by splitting cases chronologically before
186 prefix generation and by computing normalization statistics on training data only. Nonetheless,
187 residual sources of bias may persist. For example, filtering to lifecycle “complete” events standardizes
188 supervision but may discard informative “start” events that correlate with delays or cancellations;
189 the net impact on next-activity supervision is positive in our setting, yet downstream duration
190 modeling will require careful matching and robustness checks. Our compact architecture and fixed
191 hyperparameters favor reproducibility over peak accuracy; different capacity or feature sets could
192 shift the balance between Top-1 and Top-3, altering qualitative conclusions about confusion bands.

193 External validity. We evaluate on three widely used public logs that cover different control-flow and
194 resource characteristics, but they do not span the full variety of industrial settings. Domains with
195 more volatile arrivals, strict SLAs, or dynamic staffing may exhibit different ambiguity structures and
196 stronger dependence on resource policies. The simulator blueprint assumes queue observability and
197 stable activity taxonomies; in settings with concept drift, task renaming, or ad-hoc activities, both the
198 predictor and the simulator would need incremental updates and drift-aware evaluation.

199 Construct validity. Our primary metrics focus on next-activity accuracy and confusion analysis,
200 complemented by macro F1 to reflect tail classes. These are standard in PPM, yet they do not fully
201 capture operational value. For example, a model that improves Top-1 by reassigning probability mass
202 among frequent activities may have negligible effect on throughput time if the resource bottleneck
203 remains unchanged. This motivates the proposed simulator metrics—global next-event accuracy,
204 per-resource next-task precision, and workload MAPE—to better align evaluation with operational
205 objectives. Finally, we fixed Top-3 as the selection criterion for early stopping; alternative criteria such
206 as calibration error or cost-sensitive risk could yield different checkpoints with different deployment
207 trade-offs.

208 8 Conclusion

209 We provided a reproducible, leakage-safe protocol and a compact LSTM baseline for next-activity
210 prediction on three widely used event logs, together with a modular blueprint for resource-centric
211 agents implemented via per-resource policies and discrete-event simulation. The baseline delivers
212 strong Top-3 accuracy and competitive Top-1 across chronological test splits, while triangulating
213 confusion patterns highlights ambiguities consistent with unmodeled resource contention. We
214 documented pitfalls and a practical checklist spanning lifecycle handling, imbalance-aware back-offs,
215 safe indexing, and artifact logging. Next steps include releasing the full simulator with standardized
216 metrics and ablations, and integrating richer queue features and priority signals to better capture
217 operational dynamics. We hope these artifacts help the community build trustworthy, resource-aware
218 PPM experiments.

219 **References**

- 220 [1] Berti, A., van Zelst, S., and van der Aalst, W. M. P. PM4Py: Process mining for Python.
221 *SoftwareX*, 11:100110, 2020.
- 222 [2] Camargo, M., D. M. and González-Rojas, O. Learning accurate lstm models of business processes.
223 *Information Systems*, 84:1–14, 2020.
- 224 [3] Evermann, J., Rehse, J.-R., and Fettke, P. Predicting process behaviour using deep learning.
225 *Decision Support Systems*, 100:129–140, 2017.
- 226 [4] Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016.
- 227 [5] Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9(8):
228 1735–1780, 1997.
- 229 [6] Hosmer, D. W., Lemeshow, S., and Sturdivant, R. X. *Applied Logistic Regression*. Wiley, 3
230 edition, 2013.
- 231 [7] Law, A. M. *Simulation Modeling and Analysis*. McGraw-Hill, 5 edition, 2014.
- 232 [8] Tax, N., Verenich, I., La Rosa, M., and Dumas, M. Predictive business process monitoring with
233 lstm neural networks. In *CAiSE*, pp. 477–492. Springer, 2017.
- 234 [9] van der Aalst, W. M. P. *Process Mining: Data Science in Action*. Springer, 2016.

235 **Supplementary Material**

236 **A Implementation details**

237 Data loading. The discovery utility scans input directories for files with the extensions `.xes` or
238 `.xes.gz`, preferring a local input folder. XES logs are parsed with PM4Py [1] into a tidy DataFrame
239 with columns `case_id`, `activity`, `lifecycle`, `timestamp`, `resource`, with UTC timestamps and sorted by
240 time and case.

241 Prefix construction. For each case, we sort events by timestamp, filter to lifecycle “complete”,
242 and derive per-step features: inter-event delta (seconds), since case start (seconds), hour of day in
243 $[0,1]$, weekday in $[0,1]$, and a working-hours flag (Mon–Fri, 08–17). We emit prefixes of length
244 $k \in [1, \min(10, T - 1)]$ with target at position k . To avoid positional indexing bugs, timestamps are
245 converted to numpy arrays prior to computing deltas.

246 Normalization. We compute mean and standard deviation for delta and since-start on training prefixes
247 only (after chronological split), then apply the same transformation to validation and test prefixes.

248 Model and training. Activity IDs are embedded into 64 dimensions and concatenated with the five
249 continuous features, then fed to a single-layer LSTM with hidden size 128 and dropout 0.2. The final
250 hidden state goes to a linear classifier over the activity vocabulary. We train with cross-entropy loss
251 and Adam using learning rate $1e-3$, betas 0.9 and 0.999, epsilon $1e-8$, batch size 128, and select the
252 best epoch by validation Top-3 accuracy. Seeds are fixed at 42. Unless otherwise noted, there is no
253 weight decay, no label smoothing, and no gradient clipping. We use token padding with masking so
254 that loss is computed only on valid time steps.

255 Artifacts. We export per-dataset loss curves and test confusion matrices as PNGs. In the main text we
256 keep the confusion matrices and relocate loss curves to the appendix to prioritize information density;
257 additional artifact dumps are released with the code for deeper offline analysis.

258 **B Resource-centric simulator blueprint**

259 State and events. The DES maintains (i) a global clock, (ii) per-resource busy/idle status and residual
260 service times, and (iii) per-activity queues with counts and oldest waiting time. When a completion
261 event occurs, the corresponding resource becomes idle and immediately selects the next activity.

262 Policies. Each resource r has a multinomial logistic policy $\pi_r(a | x)$ over activities a with features x
 263 including previous activity executed by r , time-of-day bins, per-activity queue counts, and per-activity
 264 oldest waiting time. If samples for a class are below a threshold, we back off to a global policy π_{global}
 265 fit on all resources.

266 Durations. Each activity a has a log-normal distribution for service time with parameters fit from train-
 267 ing “start”/“complete” pairs when available; otherwise, we use the sample median from “complete”
 268 inter-event deltas as a fallback for instantaneous transitions.

269 Dispatch and ablations. Given an activity choice, the resource dispatches the oldest waiting case in
 270 the selected activity queue (FIFO within activity). The FIFO ablation replaces π_r by selecting the
 271 activity with the oldest waiting job across all queues, removing learning from the decision rule.

272 Metrics. We propose reporting: (1) global next-event accuracy comparing predicted next completion
 273 against the replayed next completion; (2) per-resource next-task precision; (3) workload MAPE
 274 comparing per-resource busy time profiles over evaluation horizons; and (4) standard PPM metrics
 275 (Top-k next-activity, remaining time MAE, suffix similarity) for completeness.

276 C Consolidated learning curves

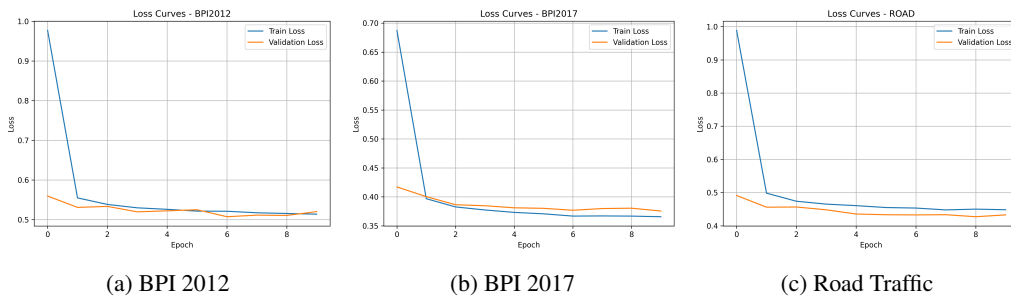


Figure 2: Training (blue) and validation (orange) loss curves for the LSTM baseline. Loss drops sharply in early epochs and then plateaus. The larger train-validation gap in BPI 2012 signals overfitting relative to BPI 2017 and Road Traffic.

277 D Broader impacts

278 This work supports reproducibility and transparency in predictive process monitoring by providing
 279 clear baselines, protocols, and simulator blueprints. The primary positive impact is lowering barriers
 280 for researchers to replicate and extend results, thereby strengthening scientific reliability. Potential
 281 risks are minimal. Overall, the societal impact is positive, contributing to more open and trustworthy
 282 research practices.

283 **Agents4Science AI Involvement Checklist**

284 1. **Hypothesis development:** Hypothesis development includes the process by which you
285 came to explore this research topic and research question. This can involve the background
286 research performed by either researchers or by AI. This can also involve whether the idea
287 was proposed by researchers or by AI.

288 Answer: [C]

289 Explanation: A postdoctoral researcher in BPM proposed the initial idea and provided a
290 short JSON note with a sketch abstract, minimal experiment outline, and key limitations. A
291 customized “AI Scientist v2” (tuned for BPM/PPM) then expanded the problem framing,
292 surveyed related work, refined the hypotheses, and generated alternative angles and ablations.
293 Human input focused on scoping and feasibility; the AI did the majority of hypothesis
294 refinement and articulation.

295 2. **Experimental design and implementation:** This category includes design of experiments
296 that are used to test the hypotheses, coding and implementation of computational methods,
297 and the execution of these experiments.

298 Answer: [D]

299 Explanation: The AI agent system produced the detailed experimental plan (data analysis,
300 splits, features, baselines/ablations, metrics, and runtime constraints) and drafted implementa-
301 tion scaffolds consistent with our BPM/PPM customization prompts. The AI contributed
302 all of the design specifics and executable structure and implemented a fully autonomous,
303 end-to-end pipeline.

304 3. **Analysis of data and interpretation of results:** This category encompasses any process to
305 organize and process data for the experiments in the paper. It also includes interpretations of
306 the results of the study.

307 Answer: [D]

308 Explanation: The AI agent system generated data analyses (tables, confusion-matrix reads,
309 error patterns, and suggested ablations) and implemented interpretation text. AI reviewed
310 for domain correctness (e.g., concurrency/resource nuances), pruned wrong statements, and
311 ensured that claims matched observed metrics and logs. Thus, AI carried the bulk of analysis
312 drafting.

313 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
314 paper form. This can involve not only writing of the main text but also figure-making,
315 improving layout of the manuscript, and formulation of narrative.

316 Answer: [D]

317 Explanation: From outline to full manuscript (sections, figures/captions, text, and references),
318 drafting was done by our AI agent system (customized version of AI Scientist v2 by
319 Sakana AI). Final polishing (clarity, tone, formatting, and minor rewrites) used ChatGPT
320 as a reviewer/editor under human supervision. Humans provided high-level guidance and
321 performed final compliance checks (style, anonymization, checklist) but did not author
322 any substantial portions of the prose nor change any claims made in the paper. This is
323 deliberately done to test the autonomy of the AI agent system for this conference.

324 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
325 lead author?

326 Description: The AI agent was able to act as a fully autonomous partner for baseline
327 construction, ablations, and reproducible experimental analysis. We deliberately avoided
328 human intervention during experiment execution and data analysis to test its autonomy.
329 While it reliably handled standard tasks and produced consistent pipelines, it struggled
330 to generate novel or complex experimental ideas beyond the templates it had been given.
331 In practice, we found it best suited as a dependable assistant for systematic evaluation
332 rather than as an originator of fundamentally new and novel methodological contributions.
333 Moreover, the references focus on well-known, established literature and the AI agent
334 system was unable to cite the latest or most directly relevant prior works. We also did not
335 add these works manually in order to not break the fully autonomous nature of the AI agent
336 system.

337 Agents4Science Paper Checklist

338 1. Claims

339 Question: Do the main claims made in the abstract and introduction accurately reflect the
340 paper's contributions and scope?

341 Answer: [Yes]

342 Justification: The abstract and introduction explicitly state that the contribution is a
343 lightweight, resource-augmented baseline. The claims are limited to reproducibility and
344 incremental performance gains, not to state-of-the-art results or proposed methodology.

345 Guidelines:

- 346 • The answer NA means that the abstract and introduction do not include the claims
347 made in the paper.
- 348 • The abstract and/or introduction should clearly state the claims made, including the
349 contributions made in the paper and important assumptions and limitations. A No or
350 NA answer to this question will not be perceived well by the reviewers.
- 351 • The claims made should match theoretical and experimental results, and reflect how
352 much the results can be expected to generalize to other settings.
- 353 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
354 are not attained by the paper.

355 2. Limitations

356 Question: Does the paper discuss the limitations of the work performed by the authors?

357 Answer: [Yes]

358 Justification: A dedicated Limitations section explains that the approach depends on re-
359 source labels, does not model concurrency or priorities, and provides only lightweight
360 proxies. It also acknowledges reduced novelty since the pipeline is primarily designed
361 as a reproducibility baseline. However, it could not implement the proposed idea fully
362 autonomously.

363 Guidelines:

- 364 • The answer NA means that the paper has no limitation while the answer No means that
365 the paper has limitations, but those are not discussed in the paper.
- 366 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 367 • The paper should point out any strong assumptions and how robust the results are to
368 violations of these assumptions (e.g., independence assumptions, noiseless settings,
369 model well-specification, asymptotic approximations only holding locally). The authors
370 should reflect on how these assumptions might be violated in practice and what the
371 implications would be.
- 372 • The authors should reflect on the scope of the claims made, e.g., if the approach was
373 only tested on a few datasets or with a few runs. In general, empirical results often
374 depend on implicit assumptions, which should be articulated.
- 375 • The authors should reflect on the factors that influence the performance of the approach.
376 For example, a facial recognition algorithm may perform poorly when image resolution
377 is low or images are taken in low lighting.
- 378 • The authors should discuss the computational efficiency of the proposed algorithms
379 and how they scale with dataset size.
- 380 • If applicable, the authors should discuss possible limitations of their approach to
381 address problems of privacy and fairness.
- 382 • While the authors might fear that complete honesty about limitations might be used by
383 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
384 limitations that aren't acknowledged in the paper. Reviewers will be specifically
385 instructed to not penalize honesty concerning limitations.

386 3. Theory assumptions and proofs

387 Question: For each theoretical result, does the paper provide the full set of assumptions and
388 a complete (and correct) proof?

389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440

Answer: [NA]

Justification: The paper does not present new theoretical results or proofs; it is an empirical baseline study.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The experiments are fully specified, including dataset selection, preprocessing, splitting strategy, feature construction, and evaluation metrics. The pipeline is deterministic and designed to be rerun reliably under time limits.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: All datasets used are standard public BPM benchmarks (e.g., BPI logs). The whole experiment pipeline is released in a single-file, self-contained form with instructions and logs for reproducing results. However, we did not want to break the fully autonomous nature of the AI agent system for this conference, so we did not add the code of AI Agent system.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Training/test splits, feature normalization, hyperparameters (e.g., logistic regression defaults), and caps (MAX_CASES, MAX_PREFIXES, label top-k) are specified in detail. This allows others to reproduce the environment.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: The study reports Top-1, Top-3, and macro-F1 metrics, but does not include error bars or statistical tests. This is because the focus is on reproducible pipeline construction under strict resource limits, not statistical comparison. Thus, the AI Agent did not implement the statistical significance tests, and we did not want to include human involvement in the results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, or overall run with given experimental conditions).

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper specifies that experiments are designed to run within a short timeout on commodity hardware (CPU, limited GPU). Explicit dataset caps and runtime constraints are included to bound compute needs.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: [Yes]

490 Justification: The work adheres to the Agents4Science Code of Ethics. It uses only public
491 datasets, avoids sensitive or private data, and does not produce outputs with foreseeable
492 negative ethical risks.

493 Guidelines:

- 494 • The answer NA means that the authors have not reviewed the Agents4Science Code of
495 Ethics.
- 496 • If the authors answer No, they should explain the special circumstances that require a
497 deviation from the Code of Ethics.

498 **10. Broader impacts**

499 Question: Does the paper discuss both potential positive societal impacts and negative
500 societal impacts of the work performed?

501 Answer: [Yes]

502 Justification: The paper discusses positive impacts (reproducible baselines, transparency in
503 BPM) in Appendix D and we discuss potential risks (overstating AI autonomy) in the AI
504 Involvement Checklist. Overall, the work supports open and trustworthy research practices.
505 Negative risks are minimal but include possible misuse of AI authorship claims without
506 transparency. These are mitigated by explicit disclosure of AI involvement.

507 Guidelines:

- 508 • The answer NA means that there is no societal impact of the work performed.
- 509 • If the authors answer NA or No, they should explain why their work has no societal
510 impact or why the paper does not address societal impact.
- 511 • Examples of negative societal impacts include potential malicious or unintended uses
512 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations,
513 privacy considerations, and security considerations.
- 514 • If there are negative societal impacts, the authors could also discuss possible mitigation
515 strategies.