

000 POSITION: WANT BETTER ML REVIEWS? STOP 001 ASKING NICELY AND START INCENTIVIZING WITH A 002 CREDIT SYSTEM 003 004

006 **Anonymous authors**

007 Paper under double-blind review

011 ABSTRACT

013 With soaring submission counts, stricter reciprocity review policies, widespread
014 adoption of platforms like OpenReview, and without the offsetting pressure of
015 publication fees, the machine learning (ML) community has one of the largest
016 scholarly presences among all scientific fields. And yet, almost *everyone* has
017 *many* unpleasant things to share about their review experience. Worse, there is
018 little public space to seriously discuss — let alone debate — what makes a review
019 system effective or how it might be improved. In this position paper, we expand
020 our discussion from the two core problems: *How can we reasonably limit the num-*
021 *ber of submissions?* and *How can we incentivize good and discourage bad review*
022 *practices?* We first assess the strengths and shortcomings of existing attempts to
023 address such problems. Specifically, we present five takes on some popular con-
024 ference mechanisms and propose two alternative designs for improvement. Our
025 general position is that meaningful improvement in ML peer review won’t come
026 from polite best-practice suggestions tucked into Calls for Papers or Reviewer
027 Guidelines — it requires **enforceable yet fine-grained procedural safeguards**
028 paired with **a currency-like credit system (what we call *OpenReview Points*)**.
029 ML practitioners can “earn” such points by contributing good review practices,
030 and “spend” across one or multiple major conferences to redeem different kinds
031 of “perks” — such as complimentary registration or the right to request additional
032 review resources.

033 1 INTRODUCTION

035 **This position paper argues that peer review in machine learning (ML) is unlikely to improve**
036 **through polite requests or optimistic guidance tucked into Calls for Papers or Reviewer Guide-**
037 **lines. Fine-grained yet enforceable procedural guardrails, combined with a spendable, across-**
038 **conference credit system, are almost mandatory for a sustainable review ecosystem.**

039 Machine learning has scaled faster than nearly any other scientific field in both volume and visi-
040 bility. We now have tens of thousands of paper submissions to a single conference,¹ open-access
041 platforms like OpenReview that support interactive discussions, and increasingly reciprocal review-
042 ing obligations to match supply with demand. On paper, the ML community has everything it needs
043 to sustain a robust yet pleasant peer-review pipeline: we have the largest scholarly presence of any
044 scientific field and the most modern review technology, all without the typical bottlenecks of pay-
045 walls, publication fees, or expensive memberships. However, the lived reality often feels far less
046 functional. From cryptic or dismissive reviews to wildly inconsistent standards, frustrations with the
047 review process are nearly universal — voiced by PhD students, seasoned professors, and industry
048 researchers alike. Worse, there is little to no structured way to hold bad actors accountable, as well
049 as incentives to encourage good actors to go the extra mile.

050 In this position paper, we expand our discussion of the two core challenges we have identified:

052 1. *How can we reasonably limit the number of submissions?*

053 ¹The most recent NeurIPS 2025 has 21,575 valid paper submissions to the main conference alone.

054
055 **2. How can we incentivize good and discourage bad review practices?**056
057 We first lay the background on why these two issues are the root causes of much unpleasantsness
058 in ML review. Then, we assess some existing attempts to mitigate such issues as implemented in
059 several ML conferences. We present our takes on such measures and, finally, propose two new
060 mechanisms: fine-grained procedural safeguards that could be enforced at scale; and a currency-
061 like incentive ecosystem based on something we called “*OpenReview Points*” — which would let
062 researchers “earn” and “spend” their reviewing efforts in tangible ways across all major conferences
063 and review cycles. We believe such mechanisms would have a fair chance of addressing many of
064 the aforementioned shortcomings effectively and, more importantly, are flexible enough to allow
065 each conference to adopt its own variants. We then present many alternative views to our proposed
066 mechanisms, where we discuss how such views are valid (or not) and how our proposed mechanism
067 shall be able to take such concerns into consideration. We conclude our paper by presenting a
068 *Recommended Practice* section, which outlines our vision on how the first few conferences adopting
069 a similar credit-based system should proceed and what aspects should be considered cautiously.
070071
072 We faithfully emphasize that our goal is not to perfect ML peer review — as it would be unfaithful
073 and condescending for anyone to claim that — but to make its failures rarer, less painful, and most
074 importantly: more accountable and sustainable.
075076 **2 ROOT CAUSES**077 **2.1 OVERBLOWN NUMBER OF SUBMISSIONS CAUSES ALL KINDS OF CHALLENGES.**078
079 We believe it is common knowledge that ML conferences typically receive an overblown amount of
080 submissions (Kim et al., 2025; Yang, 2025). Naturally, this causes all kinds of practical challenges.
081 From a manpower perspective, more submissions directly mean greater demand for reviewers and
082 Area Chairs (ACs), which translates to a heavier workload for Senior Area Chairs (SACs) and,
083 eventually, Program Chars (PCs). With such great pressure on every aspect of the conference review
084 system, the results are predictable: thinner attention per paper, more rushed triage, and greater
085 variance in both review quality and decision outcomes.
086087
088 We believe it is common knowledge that ML conferences typically receive an overblown amount of
089 submissions (Kim et al., 2025; Yang, 2025). Naturally, this causes all kinds of practical challenges.
090 From a manpower perspective, more submissions directly mean greater demand for reviewers and
091 Area Chairs (ACs), which translates to a heavier workload for Senior Area Chairs (SACs) and,
092 eventually, Program Chars (PCs). With such great pressure on every aspect of the conference review
093 system, the results are predictable: thinner attention per paper, more rushed triage, and greater
094 variance in both review quality and decision outcomes.
095 Moreover, practicality-wise, most ML conferences often guarantee the right to presentation expo-
096 sure once a paper is accepted. This makes physical capacity restrictions come into play, directly
097 imposing an upper limit on how many papers can be accepted. Many borderline or acceptance-
098 inclined papers might be ruled out purely based on capacity constraints — a role often delegated
099 to SACs. But considering the number of submissions versus the number of SACs,² this kind of
100 assignment is, by design, unreasonable and unsustainable, as no SAC has the luxury or appetite to
101 go through the content and review record of that many papers. In practice, this pressure incentivizes
102 shortcircuiting — e.g., relying more heavily on numerical scores or other similar quick heuristics —
103 which further amplifies randomness and weakens accountability. In fact, we have seen many SACs
104 publicly pushing against such “force rejection for capacity” practices, as exemplified by LinkedIn
105 posts from NeurIPS SACs Ahmad Beirami and Atlas Wang.
106107 **2.2 LACK OF OVERSIGHT, FEEDBACK LOOP, AND INCENTIVE TOWARDS GOOD AND BAD
108 ACTORS.**109
110 While reviewers, ACs, and SACs certainly have the right to provide feedback on their assigned
111 submissions, ML conferences lack proper oversight and feedback loops for such actors. A reviewer,
112 AC, or SAC can essentially engage in highly discretionary actions, so long as those actions are not
113 extreme enough to trigger a desk rejection or other formal intervention (e.g., not performing the
114 assigned review duty at all). This ecosystem leaves actors with no channel to learn how to become
115 better, let alone any real incentive to go the extra mile.
116117 As a result of the general lack of oversight, there is little to no systematic calibration across re-
118 viewers, limited visibility into meta-review quality, and minimal recognition for consistently careful
119120
121 ²Per https://media.neurips.cc/Conferences/NeurIPS2024/NeurIPS2024-Fact_Sheet.pdf, we have 195 main conference SACs but 15,671 submissions at NeurIPS 2024, making the
122 workload roughly 80 papers per SAC.
123

108 work. Conversely, low-effort or unconstructive behavior often carries no consequences; as even
 109 with escalation, there are no rules for corresponding punishment unless their action is on the most
 110 extreme end. Without routinized feedback, transparent metrics, or positive incentives, the system
 111 neither rewards exemplary stewardship nor deters poor practices. We argue that such a lack of
 112 oversight and incentive makes the overall quality drift toward the lowest-effort equilibrium, where
 113 helpful practices like internal-review discussions and thorough AC investigations rarely happen, as
 114 they are not incentivized to go those “extra miles.”

115

116 3 OUR TAKES

117

118 3.1 SOFT AND HARD SUBMISSION CAPS OFFER LIMITED HELP.

119

120 With the growing research body of the ML community (Yang, 2025; Kim et al., 2025) and with
 121 AI-assisted research becoming more accessible³ (Eger et al., 2025), the volume of submissions
 122 continues to grow at a pace that far outstrips the community’s reviewing capacity. This escalation
 123 naturally prompts discussions around mechanisms for curbing submission rates and maintaining a
 124 manageable reviewing load, where submission caps are often proposed as one of the most direct
 125 ways to reduce such volume.

126

127 We argue that submission caps — whether “soft” (e.g., mandatory reciprocal review over X submis-
 128 sions) or “hard” (e.g., strict per-author quotas on how many papers can be submitted) — provide,
 129 at best, marginal relief. The core problem is not that a small set of “hyper-prolific” lead authors are
 130 personally flooding the system with new submissions (Yang, 2025), but that there is essentially no
 131 downside to submitting unready manuscripts or endlessly recycling previously rejected work with
 132 critical flaws. We argue that, in practice, per-author caps mostly trim auxiliary authors from the by-
 133 line so that teams can fit under the quota; they do little to stop the same lead authors from submitting
 134 the same number of papers — worthy or not — or from repeatedly resubmitting flawed work. In
 135 other words, we argue that **submission caps mostly change who gets listed on a paper, rather**
 136 **than whether the paper is submitted.** This is because it would be unlikely for a team to postpone
 137 the submission of a paper simply because an auxiliary is hitting the caps. Thus, until there is a gen-
 138 uine negative incentive that discourages unlimited resubmission and rewards restraint, submission
 139 caps can only nibble at the edges of the volume problem, instead of addressing its core.

140

141 3.2 IRRESPONSIBLE REVIEWERS CARE MOST ABOUT THEIR OWN WORKS — SO ASKING 142 NICELY MIGHT NOT BE HELPFUL

143

144 The seemingly widespread presence of bad or irresponsible review practices is rooted in the lack
 145 of accountability built into current conference mechanisms. Until very recently, most ML confer-
 146 ences enforced no retaliatory punishment against irresponsible reviewers — leaving bad practices
 147 essentially unchecked. **From a reviewer’s perspective, the only thing that matters is their own**
 148 **current or future submissions. Therefore, to effectively discourage irresponsible reviewing,**
 149 **some form of penalty must be enforced at the submission end.** Otherwise, conferences have no
 150 real leverage and can only resort to asking nicely in Calls for Papers or Reviewer Guidelines; and
 151 despite the existence of some extremely thoughtful guidelines like the ARR Reviewer Guideline,⁴
 152 their effectiveness remains to be desired.

153

154 Recently, starting with CVPR 2025, several ML conferences have adopted what is essentially a
 155 retaliatory desk-rejection policy targeting irresponsible reviewers. At CVPR 2025, its Area Chairs
 156 (ACs) “*identified a number of highly irresponsible reviewers, those who either abandoned the review*
 157 *process entirely or submitted egregiously low-quality reviews, including some generated by large*
 158 *language models*” and ultimately issued desk rejections for 19 otherwise accepted papers involving
 159 those reviewers.⁵

160

161 While this act marks a meaningful start to enforcing hard procedural guardrails to protect review
 162 quality and integrity. We argue that retaliatory procedures as harsh as desk rejection can only offer

163 ³One direct piece of evidence of this might be how an AI scientist is able to get a paper accepted at the main
 164 conference track of ACL 2025. See this blog and (Zhou & Arel, 2025) for details.

165

⁴<https://aclrollingreview.org/reviewerguidelines>

⁵<https://x.com/CVPR/status/1894853624200863958>

162 marginal benefits to the conference at large, as only a few bad actors would be extreme enough to
 163 blatantly ignore direct instructions.

165 **3.3 RETALIATORY DESK REJECTION IS USEFUL, BUT IT LACKS GRANULARITY AND**
 166 **CANNOT ACHIEVE INFLUENCE AT SCALE**

168 Given the precedent set by CVPR 2025, many conferences (including ICML and NeurIPS 2025)
 169 have begun adopting similar desk rejection policies targeting ultra-irresponsible reviewers. This is a
 170 step in the right direction. But desk rejection, by nature, is a blunt instrument: it's too harsh to apply
 171 broadly and can only be reasonably used for the most extreme violations with verifiable signals. In
 172 CVPR's case, it was mostly reserved for reviewers who outright abandoned their review duties — a
 173 clean, verifiable breach that leaves no room for ambiguity.

174 Unfortunately, the vast majority of reviewer problems in ML are much more subtle than complete
 175 negligence. Irresponsible reviews can manifest in many forms: from thoughtless boilerplate com-
 176 plaints like “no theory” or “needs more experiments” applied indiscriminately to every submission,
 177 to a gross misunderstanding of basic facts and refusal to reconsider, to a raised concern with no
 178 concrete support, or even the famous “Who is Adam?”⁶ These reviews are much harder to police,
 179 but no less damaging.

180 We argue that desk rejection is too coarse a penalty to handle the long tail of poor reviewing be-
 181 haviors that fall short of full abandonment. If we want meaningful deterrents at scale, we need a
 182 system that applies graduated, proportional penalties — not just all-or-nothing rulings. The fact that
 183 the CVPR 2025 procedure only results in 19 desk rejections is clear evidence that the irresponsible
 184 review issue in the ML community is far from being resolved by merely adopting this retaliatory
 185 desk rejection policy alone; more enforceable yet fine-grained procedural safeguards are necessary
 186 to handle the wide spectrum of irresponsible review practices.

187 **3.4 100% MANDATORY RECIPROCAL REVIEWER RECRUIT IS A SLOW-ACTING POISON**

189 To keep up with rising submission counts, many ML conferences — starting with the most recent
 190 EMNLP 2025 (ARR May) — now rely on 100% reciprocal reviewer recruitment: every eligible
 191 author⁷ must also review, except in extreme circumstances like parental leave. On paper, this sounds
 192 fair: if one wants to publish, one should contribute to the review pool. But in practice, it's a slow-
 193 acting poison and done so at the cost of review quality.

194 The policy assumes that all eligible authors of every submission are both capable and willing to
 195 provide thoughtful reviews. That is simply not true under many ML literature. Many eligible authors
 196 may have played only auxiliary roles or contributed as expert consultants on highly specialized
 197 components. They are therefore ill-suited to reviewing general-purpose ML submissions. Forcing
 198 mandatory review duties on such contributors creates a predictable outcome: low-effort reviews
 199 written solely to avoid desk rejections.

200 Worse, mandatory review removes the ability for people to decline, even when they know they
 201 cannot meaningfully contribute due to sensible (but non-medical-like) emergencies. Once in the
 202 reviewer pool, conferences often allow very limited flexibility for exemption. For instance, AAAI
 203 2026 instructed their reviewers to “do your best” even if the assigned paper is outside their area of
 204 expertise. We argue that such enforced cultures would likely result in a series of rushed, templatized,
 205 or often disengaged reviews. It is worth noting that many later conferences are likely on the same
 206 page with us (in terms realizing the negative side of mandatory reciprocal review): starting from
 207 ARR July 2025, technically qualified authors may request an exemption from review duty on a
 208 case-by-case basis if they find themselves lacking the relevant expertise.⁸

209 We find it ironic that a mechanism designed to distribute the workload ends up degrading its quality.
 210 However, if everyone could opt out with no consequences, the system would collapse under the sheer
 211 volume of submissions. So, what is the reasonable middle ground — a way to allow reasonable

213 ⁶https://x.com/2prime_PKU/status/1948549824594485696

214 ⁷Where such eligibility is often determined by prior publication records, such as number of published works
 215 at A* or similar conferences.

⁸<https://aclrollingreview.org/exemptions2025>

216 opt-outs while still holding authors accountable for their share of the reviewing load? We, again,
 217 explore one such compromise in Section 4: a points-based incentive system that rewards good-faith
 218 reviewing and allows reviewers to “spend” those points to defer their reviewing obligations as one
 219 sees fit.

220

221 **3.5 HELPFUL IMPLICIT EXPECTATIONS OF ACTORS ARE OFTEN NEVER MET**

222

223 Conference processes implicitly assume that reviewers, ACs, and SACs will self-initiate best practices — such as timely calibration, substantive internal discussions, careful revision after rebuttals,
 224 and principled follow-ups — to ensure that informed decisions are made for each submission. **In**
 225 **reality, helpful practices like internal reviewer discussions almost never happen effectively be-**
 226 **cause no one is incentivized to “go the extra mile.”**

227

228 There is little real recognition, no tangible credit, and only minimal accountability for the extra
 229 coordination and time that these practices require; under deadline pressure, the rational response is
 230 to aim for the minimum viable effort. As a result, helpful practices like internal reviewer discussions
 231 become perfunctory or are entirely absent. Naturally, the decision quality then becomes noisier —
 232 not for lack of guidance, but for lack of aligned incentives to make the guidance actually happen.

233

234 **4 OUR PROPOSAL: FINE-GRAINED PROCEDURAL GUARDRAILS WITH A**
 235 **CURRENCY-LIKE INCENTIVE SYSTEM**

236

237 To make meaningful progress in peer review reform, we argue that two ingredients are essential:
 238 **enforceable procedural safeguards at different granularities**, and **an incentive structure that**
 239 **rewards good-faith participation while offering flexibility**. We propose a system based on a
 240 community-wide, cross-conference-supported economy called “OpenReview Points” — mainly for
 241 the mainstreamness of OpenReview and its good position to keep track of such balance.

242

243 This section outlines the basic principles of such a system, discusses potential enforcement strate-
 244 gies, and explores the feasibility of a conference-wide credit market that could finally provide con-
 245 ferences organizers both the “stick” and the “carrot” they currently lack.

246

247 **4.1 OPENREVIEW POINTS: A CURRENCY-LIKE ECONOMY ENABLING FLEXIBLE OPTIONS**

248 The current review ecosystem operates on the honor system — a reviewer is expected to perform
 249 review duties diligently and hope that others will do so as well. However, we argue that optimistic
 250 hoping is not a system. To install accountability, we propose currency-like credit system, giving
 251 contributors to the review pipeline something to earn, spend, and track.

252

253 Under our proposal, ML practitioners would accumulate OpenReview Points based on their contribu-
 254 tions to the community. For instance — in the context of reviewing — completing a standard
 255 review might earn 1 point, helping with an emergency review might earn 2, and being recognized
 256 as an “outstanding reviewer” could grant an additional 3 points.⁹ Once earned, OpenReview Points
 257 could be spent to gain access to certain “perks” and privileges. For example:

258

- A reviewer can spend 5 points to opt out of an assigned review duty.
- An author can spend 10 points to exempt a co-author from their reciprocal reviewing obligation.
- An author can spend 50 points to request an additional expert reviewer in the case of a highly controversial or borderline decision. In the meantime, a reviewer/AC can take this job and earn those 50 points.
- An author can spend 100 points to redeem free registration.

259

260

261

262

263

264

265

266

267

268

269

263 This economy introduces direct incentives: if one contributes meaningfully, one gains flexibility and
 264 optionality. If one does not, one’s publishing privileges will begin to shrink. It also gives conference
 265 panels more space to experiment with different policies and enforcement harshness, without relying
 266 on blunt-force policies like universal reciprocal reviewing or desk rejections.

267
 268 ⁹We emphasize that all point values mentioned in this section are intuitively assigned for hypothetical pur-
 269 poses. A real OpenReview Point-based economy would require significantly more sophisticated balancing,
 subject to each conference’s own preferences. More on such specifications in Section 6.

270 For instance, much of our work argues that there is a lack of incentive for actors to “go the extra
 271 mile,” even when such effort can be immensely helpful (e.g., as anecdotally demonstrated in Ap-
 272 pendix C). With point incentives, however, such “extra miles” can be encouraged: reviewers shall
 273 become more willing to initiate and engage in internal discussion, and ACs shall become more will-
 274 ing to investigate — simply because exemplary actions can now be rewarded. We can push this
 275 further by introducing targeted awards and penalties to shape community behavior. As discussed in
 276 Section 2.2, the lack of a reviewer feedback loop can potentially be mitigated by awarding points to
 277 authors who provide detailed reviewer feedback that reviewers may consult to improve their future
 278 practices. Similarly, as noted in Section 2.1, many reviewing issues stem from inflated submission
 279 counts. **One way to mitigate this is to require a small and refundable “submission fee” — e.g.,**
 280 **10 OpenReview Points — per paper.** If the paper is accepted or meets a reasonable “fair attempt”
 281 bar, the points are refunded; otherwise, they are forfeited. This soft deterrent discourages unready
 282 submissions by linking low-quality or premature work to a corresponding reduction in future publi-
 283 cation privileges.

284 **To be clear, we are not arguing for the enforcement of any specific rule** — whether that is
 285 charging a submission fee in points, or allowing opt-outs from reviewing, or more. Rather, we argue
 286 that a currency-like system would grant every participant in the ML community far greater flexibility
 287 in how they interact with the review process. **While we fully expect friction or disagreement**
 288 **regarding any particular rule or redemption policy, we believe it would be difficult to argue**
 289 **against the utility of having a credit-based system at all.** Since it makes sense for different
 290 conferences to carve out their own rules to cater to their own communities.

291 4.2 VOTING-BASED PENALTIES MAY INTRODUCE FALSE POSITIVES — AND THAT’S 292 ACCEPTABLE, BECAUSE WE GET FINE-GRAINED ENFORCEMENTS IN RETURN 293

294 To discourage low-effort or malicious reviewing, we propose allowing area chairs and fellow re-
 295 viewers to flag irresponsible reviewer behavior with much greater flexibility and finer granularity.
 296 Of course, actions that are verifiably egregious — such as completely missing reviews or openly
 297 posting LLM-generated content — should trigger the harshest penalties (e.g., desk rejection), as
 298 these cases are binary, easy to verify, and largely undisputed.

299 **But most bad reviewing practices do not look like that. They are far more subtle:** templated
 300 one-liners like “no theory” or “needs more experiments” applied indiscriminately, vague dismissals
 301 without justification, raising a plethora of out-of-scope questions, or reviews that clearly misunder-
 302 stand the paper’s core contributions and never bother to revise... **These cases are harder to catch**
 303 **algorithmically and rarely rise to the level where desk rejection is justifiable. Yet, they are**
 304 **widespread and deeply harmful.** As authors are often forced to invest time and energy to address
 305 these reviews, which is often to little effect if no oversights are placed upon the reviewers’ end.

306 This is where a voting-based penalty system can help. For instance, if the authors report a reviewer,
 307 and that reviewer’s peers on the paper — along with the area chair — unanimously agree that a
 308 review is deemed unacceptably low in quality, that reviewer could receive penalties ranging from a
 309 warning to various levels of point deduction.

310 We recognize that such systems introduce the possibility of false positives. But in the meantime, we
 311 argue that it is largely acceptable: First, if safeguards such as unanimous agreement, AC confirma-
 312 tion, and an appeal mechanism are in place, we believe the practical false positive rate can be kept
 313 low. Secondly, since point deduction is a much less extreme act compared to penalties like desk
 314 rejection, even a false case is unlikely to result in severe and immediate impacts of irrecoverable
 315 harm. More importantly, while our proposed system is far from perfect, the current system basically
 316 has the opposite problem: it has a nearly 0% true positive rate — since no matter how badly the
 317 reviewer behaves, as long as it is not at an automatically verifiable level of atrociousness, there are
 318 zero consequences. The status quo is clearly worse in a comparative sense.

319 No penalty system will ever be perfect, but the absence of one guarantees stagnation. We argue
 320 that a small risk of overcorrection is a worthwhile price to pay for finally holding the peer-review
 321 process to a higher standard. An even lower risk alternative is to award credits to AC/reviewers for
 322 providing detailed feedback on their peer reviews. this would enable a positive feedback loop where
 323 actors have channels to learn and become better versions of themselves. We discuss more about such
 specific recommended practices and our considerations behind them in Section 6.

324

5 ALTERNATIVE VIEWS

325
326 While we advocate for enforceable procedural safeguards and a currency-like incentive system, we
327 recognize that not everyone will agree with this approach. Below, we discuss several alternative
328 perspectives and respond to their concerns.
329330

5.1 “PEER REVIEW SHOULDN’T BE GAMIFIED.”

331 A common objection is that introducing a credit system risks gamifying the review process — turning
332 what should be a scholarly, community-driven responsibility into a transactional system. While
333 we sympathize with this concern, our counterpoint is simple: peer review is already governed by
334 incentives — such as reviewing others’ submissions in turn for having one’s own work properly
335 reviewed — but these incentives are just poorly aligned and implicitly stated.
336337 Researchers submit to conference because they care deeply about getting their own work accepted;
338 yet, they are often disincentivized from reviewing carefully and consistently. A credit system does
339 not create incentives out of thin air — it simply formalizes them and aligns them with the broader
340 health of the ecosystem.
341342

5.2 “VOTING-BASED PENALTIES WILL BE ABUSED OR POLITICIZED”

343 Another concern is that a flagging or voting-based penalty system could be misused — weaponized
344 in borderline cases or influenced by interpersonal bias. We agree that any enforcement mechanism
345 needs guardrails. That’s why we require close-unanimous (if not totally unanimous) agreement from
346 all other participating reviewers on the same paper and area chair confirmation before any penalty is
347 issued. False positives are not impossible, but they are rare and correctable when hedged with these
348 guardrails, and their damage — some point deduction — is unlikely to cause irrecoverable harm like
349 immediate desk rejections and submission bans. The alternative — a system that allows actors to
350 act with no accountability whatsoever — is far more damaging in the long run.
351352

5.3 “A POINT SYSTEM FAVORS THE PRIVILEGED ACTORS, AND CAN RESULT IN BAD THINGS 353 LIKE LOSING QUALITY REVIEWERS.”

354 Some may argue that a credit system will disproportionately benefit researchers with more time,
355 institutional support, or prior connections — allowing them to “buy” their way out of responsibilities
356 (e.g., being exempt from review duties) while leaving others to “pick up the slack.” This is a fair
357 concern. But in our design, points are earned through labor, not status. There is no “premium tier
358 of citizen,” only accumulated contributions through hard work. While it is still true that researchers
359 with strong support will likely have more opportunities to contribute — as they are not otherwise
360 occupied by some chores — their “surplus contributions” are still a net gain to the community.
361362 We also note that, while exemptions from review duties might indeed result in losing reviewers, one
363 thing to consider is whether those who are willing to pay a high price to be exempted are producing
364 quality reviews (if kept by force), and whether they are likely to have the bandwidth to stay engaged
365 with the authors. We argue that a better alternative might be to just let them be exempted, and
366 utilize the collected points to incentivize reviewers who do have the bandwidth and motivation in
367 this particular conference cycle.
368369 On the same note, one thing we would strongly advocate is **mobilizing researchers who are not**
370 **main authors to participate more in the review pipeline**, as they likely have better bandwidth
371 (since they are not under the pressure of author deadlines), and their reviews will not be as affected
372 by feedback on their own submitted work. Under the current system, there is little incentive for
373 researchers to do so, as most reviewers are recruited by mandatory reciprocity, which no longer
374 applies without being an author. Our proposed credit system might provide them with a strong
375 incentive to participate, as they can earn points to enrich their publication privileges; and specifically,
376 have the option to spend such points to be exempt from reviewer duties when they are submitting
377 lead-authored work, granting themselves wider bandwidth as authors when they are under rebuttal
pressure.
378

378 5.4 “ALL OF THIS SOUNDS TOO BUREAUCRATIC.”
379380 Some may worry that procedural enforcement, tracking points, and adjudicating review quality will
381 introduce too much bureaucracy into the process. We, again, see where this concern is coming
382 from. However, conferences already invest massive effort coordinating thousands of reviews and
383 rebuttals; we are simply proposing mechanisms to make those efforts fairer, more consistent, and
384 more sustainable over time. While we do agree that the full form of our credit system can be too
385 heavy to be implemented at once, we argue a gradual roll out of changes can be rather “soft landing”
386 to existing community members.
387388 5.5 “A CROSS-CONFERENCE RECIPROCITY MUST EXIST FIRST”
389390 One clear and fair criticism of our credit system is that, for it to work to its full potential, multiple ma-
391 jor conferences must adopt it. Granted that conferences like ICML, NeurIPS, and ICLR almost never
392 work together closely, we recognize that such a prerequisite can be difficult to achieve. However, we
393 argue that there are ML conferences well-positioned to adopt such practices: e.g., the ARR series of
394 conferences has long implemented cross-conference measures (e.g., submission bans from the next
395 ARR cycle), as experimented with in EMNLP 2025,¹⁰ making them more openminded to adopting
396 other similar cross-conference measures. Further, even if the credit system is per-conference, it can
397 still function at a level that is better than nothing; it is just that features requiring accumulated effort
398 may be harder to activate and experiment with.
399400 5.6 “WHAT ABOUT EARLY-CAREER RESEARCHERS?”
401402 Much like how companies and games handle the onboarding process for novice players and new
403 employees, a reasonable expectation is that the point-hosting platform shall grant a baseline amount
404 of credits to such first-time contributors, perhaps also with a protection period (much like how new
405 hires cannot be directly placed into a Focus/PIP pipeline in tech companies), so that they have
406 enough time and capital for trial-and-error.
407408 Seeking endorsement from a seasoned scholar (much like how arXiv operates) is another viable
409 option, though it might be more unfriendly to researchers from underrepresented backgrounds. For
410 such contributors, a platform might run routine workshops, where participating and passing such
411 workshop education shall grant these contributors such baseline credits.
412413 6 RECOMMENDED PRACTICES
414415 As emphasized throughout this position paper, our goal is not to promote a single, prescriptive
416 rulebook that every conference must follow, but to advocate for a flexible framework that can adapt
417 to different conference idiosyncrasies. Every policy comes with its own trade-offs — it is, ultimately,
418 an art of compromise to decide which set of policies to adopt. **Under our credit system, conference**
419 **panels and authors are akin to store owners and customers: the panels decide what goods are**
420 **offered and at what price, while the customers decide where and how they wish to spend their**
421 **money.**
422423 However, we fully recognize that without concrete discussion of how such a system might actually
424 operate, there will naturally be resistance — and, worse, chaos if adopted without proper considera-
425 tion. Thus, this section serves as practical guidance from us authors on how the first few conferences
426 adopting a credit-like system might proceed. We also leave our answers to many frequently asked
427 questions to Appendix B
428429 6.1 HEAVY ON EXISTING PERKS
430431 We believe that early adopters should anchor their credit-like system around existing perks already
432 offered in current ML conferences (e.g., complimentary registration, emergency reviewer invita-
433 tions), rather than immediately introducing entirely new perks. Staying with existing perks provides
434¹⁰<https://2025.emnlp.org/reviewer-policies/>

432 two immediate benefits: 1) Because these perks are already part of established workflows, redistributing them according to the point-based system (e.g., ranking reviewers by awarded points per
 433 conference cycle rather than relying on an AC’s subjective judgment) keeps overall impact bounded.
 434 If the new distribution turns out problematic, its effects are still confined to the known scope of
 435 these already-tested perks; whereas brand-new perks introduce unknown risks. 2) We can directly
 436 compare key metrics between credit-system conferences and their historical data. Any improve-
 437 ment or decline is then more likely attributable to the credit system (or its specific implementation),
 438 rather than being confounded by the introduction of new perks. Using the same set of perks but with
 439 credit-based allocation gives us an opportunity to collect data on how the credit system behaves in
 440 practice. Such data shall serve as the baseline to test out different new perks.
 441

442 443 6.2 GENTLE AND GRADUAL ROLLOUT OF NEW PERKS, POTENTIALLY WITH ONE-OFF TESTS

444 When launching new perks, it is best to roll them out gently and gradually rather than all at once.
 445 This approach offers a clean testbed to monitor each perk’s contribution and reduces the information
 446 load on all involved parties, who will need time to adjust.
 447

448 Observant readers may notice that some of our proposed policies — e.g., free registration, the right
 449 to request additional reviewers, or refundable submission fees discussed in Section 4 — require a
 450 relatively long-term accumulation of points to become useful. Researchers would need to earn (or be
 451 penalized) enough points before crossing the thresholds of these services, stretching the evaluation
 452 horizon.

453 A simple way to expedite early evaluation is to introduce one-off tests. For example, in addition
 454 to point awards, conferences might grant top point-earners a one-off right to request an additional
 455 reviewer to help resolve borderline cases, with this right expiring at the end of the conference cycle.
 456 This allows organizers to directly observe whether such redeemable incentives meaningfully help —
 457 and to what extent these improvements propagate through the reviewer–AC pipeline. This kind of
 458 fast feedback might help conference organizers trim unhelpful policies quickly and expedite a faster
 459 iterations of rule sets.
 460

461 462 6.3 DETERMINING POINT VALUES FOR CONTRIBUTIONS AND PERKS

463 One reason we did not specify exact point values for different contributions (e.g., how many points
 464 an emergency review should yield) is that we currently lack the empirical data needed to set these
 465 responsibly — such as the typical number of emergency review requests in a conference, how often
 466 these requests are fulfilled, and how their workload compares to regular reviews. As a general
 467 guideline, however, we believe it is reasonable to treat the completion of one regular review duty as
 468 the base “unit price” of this ecosystem.

469 For the sake of discussion, suppose one regular review yields 1 OpenReview Point. Rather than
 470 arguing directly about how many points each contribution “should” receive, we propose working
 471 backwards from the value of the *perks* the system must support. For example, we can roughly
 472 estimate the point cost of a free registration by (1) checking how many free registrations a conference
 473 can realistically offer,¹¹ (2) determining what fraction of the reviewer population this corresponds
 474 to (e.g., the top 5–10%), and (3) computing the total point budget issued in that conference cycle —
 475 roughly, the number of submissions times the number of required reviews, plus bonus points (e.g.,
 476 awards) and points collected via exemption fees.

477 With this budget in hand, we can set the point cost of a free registration so that only the intended
 478 contributor percentile can redeem it under plausible participation patterns.

479 Once the point cost of a perk is anchored, the point values of contributions should be calibrated
 480 relative to it. For instance, if completing one’s regular review duty plus two emergency reviews
 481 typically places a reviewer in the top percentile eligible for free registration, then the point-value of
 482 an emergency review should be set accordingly. While this may still leave questions like why two
 483 emergency reviews, rather than three or five, should map to this threshold, working from anchored
 484 perks at least constrains point allocations to a meaningful, narrow range.

485 ¹¹For instance, NeurIPS 2025 offers free registration to top reviewers, which amounts to around 1,900+
 486 reviewers. Similar travel grants and volunteer opportunities are also offered at most ML conferences.

486 In other words, point-value determination begins by anchoring perk values to the scale of the econ-
 487 omy; contribution-level point allocations then follow to maintain internal consistency rather than
 488 being chosen in isolation.

490 **6.4 TRACK KEY METRICS AND PUBLICIZE SUCH STATISTICS**

492 Finally, for a credit system to have a lasting impact, conferences must make informed decisions
 493 about which rules to adopt and at what point-values. Such decisions require cross-conference con-
 494 sistency. If, for example, NeurIPS values its perks at 10x ICLR’s level for no meaningful reason, the
 495 ecosystem loses the interoperability we envision. Thus, each conference should monitor key metrics
 496 and publish these statistics as part of their post-conference fact sheets.

497 For example: If a new rule is implemented, do we observe increased interaction among reviewers
 498 and ACs? Do ACs report that these additional exchanges help them make more confident decisions?
 499 These statistics and reports will form the foundation for iterating toward a better implementation of
 500 the credit system and will also serve as a strong signal — even an advertisement — encouraging
 501 more conferences to adopt a shared credit currency.

502 **7 CONCLUSION AND LIMITATIONS**

505 In summary, we present a flexible credit-based framework aimed at improving accountability,
 506 aligned incentives, and procedural fairness in the ML peer-review pipeline. While we believe the
 507 system offers practical value and a principled foundation for future experimentation, it is equally
 508 important to acknowledge its current limitations and the challenges in evaluating such proposals.

510 **Lack of numerical support and simulated experiments** Our work lacks numerical results, but
 511 we believe discussions about review mechanisms are most meaningful in the hypothetical space —
 512 since there is no way to rewind history and A/B-test two parallel conference processes. Likewise,
 513 LLM-powered simulations add little value: too many layers of prompt, decoding, and model-specific
 514 variance compound, making such outcomes highly sensitive and difficult to trust. For instance,
 515 even if we scraped full ICLR discussion logs, what would we meaningfully do with them?

- 516 • Prompt LLM “reviewers” to apply our credit system and judge their reviews with another LLM?
- 517 • Ask roleplaying LLMs to simulate reviewer discussions and treat their frequency as a signal?
- 518 • Simulate multi-conference point accumulation and claim reviewer-load reduction?

519 We believe most readers would view such simulations as overly fragile and of limited interpretabil-
 520 ity. A further complication is that many key metrics are not public. Useful signals such as the
 521 number of internal reviewer discussions or reviewer-profile dynamics are unavailable. Even crude
 522 simulations — e.g., prompting LLMs to decide whether submissions fall below a “fair attempt bar”
 523 and deducting points — quickly become hand-wavy; as changing a threshold or banning simulated
 524 authors from a “next conference” would trivially reduce load, but such results would be meaningless
 525 without real baselines.

526 That said, we recognize the desire for some anchoring to real conferences. **Thus, we share three**
 527 **case studies — drawn from similar top ML venues where we served as reviewers — in Ap-**
 528 **pendix C.** While works without statistical evaluation may face a higher bar for ICLR main track,¹²
 529 we believe our contribution fits ICLR’s call for “*how we can improve the ways that we conduct*
 530 *and evaluate machine learning research*” and offers tangible value to the community — even if it
 531 improves the review pipeline only modestly.

532 **Position paper submitted to ICLR main track** Although this paper is written in a typical
 533 position-paper style, ICLR does not have a dedicated track for such works. We have confirmed
 534 with the PCs that “*position papers can be published at ICLR if they have sufficient novelty and*
 535 *value for the ICLR community*,” as demonstrated by prior accepted work (Ngo et al., 2024). We
 536 argue that ICLR is the most appropriate venue among the ICML/NeurIPS/ICLR trifecta — both in
 537 progressiveness and in its community’s willingness to engage with systemic issues that affect all
 538 researchers. **We ask for reviewers’ open-mindedness and invite active discussion.**

539 ¹²<https://openreview.net/forum?id=fh8EYKFKn> is one accepted example.

540
541 **ETHICS STATEMENT**542
543 Granted our work operates in a hypothetical space and promotes no particular ML methods or
544 datasets, it is our honest analysis that it requires no further ethical consideration. We have care-
545 fully read the ICLR Code of Ethics and pledge to adhere to it.546
547 **THE USE OF LARGE LANGUAGE MODELS (LLMs)**548
549 We would like to disclose that part of the writing of this paper was polished by a language model,
550 though a human researcher is there to verify that the final output is true to the researcher's opinion.551
552 **REFERENCES**553
554 Alina Beygelzimer, Yann N Dauphin, Percy Liang, and Jennifer Wortman Vaughan. Has the ma-
555 chine learning review process become more arbitrary as the field has grown? the neurips 2021
556 consistency experiment. *arXiv preprint arXiv:2306.03262*, 2023.557
558 Corinna Cortes and Neil D Lawrence. Inconsistency in conference peer review: Revisiting the 2014
559 neurips experiment. *arXiv preprint arXiv:2109.09774*, 2021.560 Steffen Eger, Yong Cao, Jennifer D'Souza, Andreas Geiger, Christian Greisinger, Stephanie Gross,
561 Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, et al. Transforming science with large lan-
562 guage models: A survey on ai-assisted scientific discovery, experimentation, content generation,
563 and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.564
565 Armen Yuri Gasparyan, Alexey N Gerasimov, Alexander A Voronov, and George D Kitas. Re-
566 warding peer reviewers: maintaining the integrity of science communication. *Journal of Korean*
567 *medical science*, 30(4):360–364, 2015.568 Alexander Goldberg, Ivan Stelmakh, Kyunghyun Cho, Alice Oh, Alekh Agarwal, Danielle Bel-
569 grave, and Nihar B Shah. Peer reviews of peer reviews: A randomized controlled trial and other
570 experiments. *PloS one*, 20(4):e0320444, 2025.571
572 Steven Jecmen, Nihar B Shah, Fei Fang, and Leman Akoglu. On the detection of reviewer-author
573 collusion rings from paper bidding. *Transactions on Machine Learning Research*, 2025. ISSN
574 2835-8856.575
576 Jaeho Kim, Yunseok Lee, and Seulkil Lee. Position: The AI conference peer review crisis de-
577 mands author feedback and reviewer rewards. In *Forty-second International Conference on Ma-*
578 *chine Learning Position Paper Track*, 2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=18QemUZaIA)
579 [id=18QemUZaIA](https://openreview.net/forum?id=18QemUZaIA).580
581 Ryan Liu and Nihar B Shah. Reviewergpt? an exploratory study on using large language models for
582 paper reviewing. *arXiv preprint arXiv:2306.00622*, 2023.583
584 David Mimno and Andrew McCallum. Expertise modeling for matching papers with reviewers. In
585 *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and*
586 *Data Mining*, KDD '07, pp. 500–509, New York, NY, USA, 2007. Association for Computing
587 Machinery. ISBN 9781595936097. doi: 10.1145/1281192.1281247. URL <https://doi.org/10.1145/1281192.1281247>.588
589 Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning
590 perspective. In *The Twelfth International Conference on Learning Representations*, 2024. URL
591 <https://openreview.net/forum?id=fh8EYKFKns>.592
593 Vishisht Rao, Justin Payan, Andrew McCallum, and Nihar B. Shah. ML researchers support openness
594 in peer review but are concerned about resubmission bias, 2025. URL <https://arxiv.org/abs/2511.23439>.

594 Anna Rogers and Isabelle Augenstein. What can we do to improve peer review in NLP? In Trevor
 595 Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1256–1262, Online, November 2020. Association for Computational
 596 Linguistics. doi: 10.18653/v1/2020.findings-emnlp.112. URL <https://aclanthology.org/2020.findings-emnlp.112/>.

597

598 Buxin Su, Jiayao Zhang, Natalie Collina, Yuling Yan, Didong Li, Kyunghyun Cho, Jianqing
 599 Fan, Aaron Roth, and Weijie Su. The icml 2023 ranking experiment: Examining author self-
 600 assessment in ml/ai peer review. *Journal of the American Statistical Association*, 0(0):1–12, 2025.
 601 doi: 10.1080/01621459.2025.2510006. URL <https://doi.org/10.1080/01621459.2025.2510006>.

602

603 Weijie J Su. You are the best reviewer of your own papers: An owner-assisted scoring mecha-
 604 nism. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in
 605 Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=xmx5rE9QP7R>.

606

607 Jibang Wu, Haifeng Xu, Yifan Guo, and Weijie Su. A truth serum for eliciting self-evaluations in
 608 scientific reviews. *arXiv preprint arXiv:2306.11154*, 2023.

609

610 Jing Yang. Position: The artificial intelligence and machine learning community should adopt a
 611 more transparent and regulated peer review process. In *Forty-second International Conference on
 612 Machine Learning Position Paper Track*, 2025. URL <https://openreview.net/forum?id=gnyqRarPzW>.

613

614 Yichi Zhang, Fang-Yi Yu, Grant Schoenebeck, and David Kempe. A system-level analysis of confer-
 615 ence peer review. In *Proceedings of the 23rd ACM Conference on Economics and Computation*,
 616 pp. 1041–1080, 2022.

617

618 Andy Zhou and Ron Arel. Tempest: Autonomous multi-turn jailbreaking of large language models
 619 with tree search. *arXiv preprint arXiv:2503.10619*, 2025.

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648 A RELATED WORKS 649

650 **Proposal of conference mechanisms for better ML review quality** To the best of our knowledge,
 651 few published works have addressed **how to improve ML review quality by adjusting conference**
 652 **mechanisms**. The closest work to ours might be Kim et al. (2025), where the authors advocate
 653 establishing a feedback loop to reviewers and promoting reviewer rewards — two points that we
 654 also share. Specifically, Kim et al. (2025) highlights that if we allow authors to rate reviewers, those
 655 ratings will almost always be heavily influenced by the specific strengths and weaknesses (and, by
 656 extension, the scores) listed by the reviewer (Goldberg et al., 2025). Reviewers who fairly rate
 657 papers negatively may be subjected to unfair retaliatory ratings from authors. To address this, Kim
 658 et al. (2025) suggests a two-stage reveal: authors first read only the reviewer-written summary and
 659 strengths and provide a rating, where the weaknesses are then visible. We argue that this system
 660 might work to some degree, but in reality, much of a review’s quality is determined by whether the
 661 highlighted weaknesses are sound and well supported. Reviewing reviewers without seeing these
 662 details would likely produce a lot of noise, and reviewers would be incentivized to write vague
 663 summaries that lack sharp substance. As for reviewer rewards, Kim et al. (2025) mostly argues
 664 for vanity perks like digital badges (e.g., ones similar to the “Pull Shark” badge on GitHub). We
 665 argue that such vanity-only perks would have much less influence than the review, submission, and
 666 cost-influencing perks we propose here.

667 Another piece of related work is the Isotonic Mechanism score pioneered by Su (2021) and its
 668 follow-up works like Wu et al. (2023); Su et al. (2025), which survey authors (with multiple submis-
 669 sions) and ask them to rank their submitted papers. A score is thereby calculated and compared with
 670 the mean of reviewers’ raw scores. Should these two scores exhibit too drastic a gap, it may trigger
 671 AC intervention with an additional reviewer request, etc. We note that this work is, by large, orthog-
 672 onal to ours, as it is yet another safeguard one can implement under our proposed credit system. The
 673 two works overlap in the sense that some countermeasures do show resemblance (e.g., requesting
 674 an additional reviewer).

675 Two pieces of relatively early but directly related work are Rogers & Augenstein (2020) and Zhang
 676 et al. (2022). Rogers & Augenstein (2020) discusses why incentives clash under a peer-review con-
 677 text and outlines many potential proposals (e.g., better review–paper matching, more tracks, abol-
 678 ishing “score-based” feedback, track-specific review formats, etc.). Similarly, Zhang et al. (2022)
 679 also investigates various policies but focuses on modeling why resubmission is so prevalent despite
 680 many works eventually getting accepted at a top venue. The main difference between these works
 681 and ours is that they do not propose a unified recipe (e.g., a credit system) as a means to address
 682 various issues; instead, custom solutions are typically proposed and discussed for each problem or
 683 edge case.

684 **Conference Review Statistics** There are a few works like Yang (2025); Cortes & Lawrence
 685 (2021); Beygelzimer et al. (2023); Goldberg et al. (2025) that collect real statistics and conduct con-
 686 trolled experiments from past conferences. While such works typically do not propose mechanism-
 687 based solutions, their numerical presentations help illustrate the scale and muddiness of current
 688 conferences; yet, such controlled exploration shall offer us insight into the practical dynamics of a
 689 particular mechanism design.

690 **Boarder Relevant Art** Last, outside the machine learning community, we have Gasparyan et al.
 691 (2015) analyzing peer review incentives under a mainly medical-focused context. The main argu-
 692 ment of this work is *“none of these (financial or nonfinancial incentives) is proven effective on its*
 693 *own”*; however, the authors envisioned *“a strategy of combined rewards and credits for the review-
 694 ers’ creative contributions seems a workable solution.”*

695 Our work makes essentially the same argument, though framed under a point-based system (and of
 696 course, with more ML-specific flavors). Our system supports many of the typical financial and non-
 697 financial incentives mentioned in Gasparyan et al. (2015). For instance, several exemplary policies
 698 we discuss in Section 4.1 range from incentives for reviewers to write higher-quality and more timely
 699 reviews, to deterrents for authors submitting unready work, to non-financial privileges such as the
 700 right to be exempt from review assignments, and even financial compensation like free registration.
 701 We believe it is fair to say that we are not advocating for any particular type of incentive, but rather

702 a collection of them. It just so happens that we unify them under a credit-based system, allowing
703 their impact to last beyond a single conference.
704

705 Another point specific to Gasparyan et al. (2015) is that many reviewers dismiss incentives such as
706 paper purchase discounts or free publication access, since their institutions already pay for publisher
707 subscriptions, making such incentives mostly relevant to members with non-academic backgrounds.
708 In our proposal, however, conference hosts can offer services that no institution would purchase in
709 bulk (e.g., registration), or even privileges that money cannot buy (e.g., the right to request additional
710 review resources). While we do not claim these incentives are inherently better — as one man’s
711 vulgarity is another’s lyric — we do believe our framework offers a flexible way to combine different
712 incentives to suit diverse needs, effectively pushing forward a system that Gasparyan et al. (2015)
713 envisioned.

714 Outside the effectiveness of a certain incentive system, much of a conference’s experience also
715 depends on many other aspects, such as the reduction of collusion rings (Jecmen et al., 2025), finding
716 better reviewer–paper mappings (Mimno & McCallum, 2007), determining the level of openness
717 (Rao et al., 2025), or the use of LLMs for reviewers (Liu & Shah, 2023). We refer users to such
718 works to build a deeper understanding of these important topics. We recommend readers refer to
719 Kim et al. (2025) for an overview of such works.
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 **B FREQUENTLY ASKED QUESTIONS**
757758 **Who governs this?** We believe it is best to have a bank-like entity to govern point storage. Neutral
759 platforms like OpenReview could serve as an ideal medium for tracking how many points each
760 author has accumulated.
761762 **How do we prevent bad behavior (fraud, point farming, unfair allocation, etc.)?** No mecha-
763 nism can fully eliminate misconduct. Even real-world economies with tangible consequences face
764 persistent bad actors. It would be disingenuous for us — or anyone — to claim otherwise. However,
765 any thriving economy depends on aligning the interests of its participants. Conference organiz-
766 ers and authors share the goal of obtaining high-quality reviews and meaningful discussion, while
767 reviewers (in the worst case) may simply want to earn points.
768769 To bridge these goals, additional points can be awarded for high-quality reviews, and penalties
770 applied for low-quality ones. We can also limit the total review capacity per author to discourage
771 “quick review for basic point farming.” Enforcement can be driven by an AC + reviewer voting
772 mechanism, as discussed in various parts of our paper. We believe that these enforceable safeguards
773 would be effective in mitigating typical point-farming behaviors where review quality is severely
774 compromised.
775776 **Would Y operation be allowed (e.g., point transfer)?** Point transfers should be rare and limited,
777 as the foundation of this system is that contributors must personally engage in community service
778 to earn points. However, some team-based privileges make sense. For instance, authors of the same
779 paper could collectively “purchase” certain “products” — e.g., request additional review resources
780 for borderline cases, or exempt a non-expert coauthor from review duties.
781782 In line with our goal of proposing a framework rather than a detailed rulebook, we find it reasonable
783 to forbid person-to-person transfers while allowing team-based purchases or awards, depending on
784 conference preferences.
785786 **Would Z infra/consensus be needed?** It is fair to note that certain infrastructure and shared con-
787 sensus would be required to operationalize our credit system. At minimum, platforms like OpenRe-
788 view should provide:
789790

- A balance-tracking system to record point accumulation per author.
- A limited point-transfer interface to implement awards or penalties.

791792 While such infrastructure is necessary, we emphasize that it is a minimal requirement. The “legisla-
793 tion” and enforcement of specific policies would still rest with individual conference panels. As one
794 man’s vulgarity is another’s lyric, there will never be a perfect policy — only carefully considered
795 trade-offs. Our framework respects the diverse needs of the community and supports them with
796 great flexibility.
797798
799
800
801
802
803
804
805
806
807
808
809

810 C THREE CASE STUDIES WHERE WE SERVE AS THE REVIEWERS
811
812813 While our position paper, unfortunately, lacks real conference data to support why our proposed
814 framework would be helpful, we believe there is even less reason to run an LLM-roleplaying sim-
815 ulation. We understand the perspective that having some anchors to real conferences is preferred.
816 Here, we share three case studies — all from similar top ML conferences — where we serve as
817 reviewers.
818
819820 C.1 CASE 1: REVIEWERS CRITICIZING MATTERS OUTSIDE THE PAPER'S SCOPE.
821
822823 In this case, we observed that another reviewer criticized the submitted work for reasons clearly
824 outside its intended scope. We therefore raised our concern to the AC:
825
826

Internal comment to PC/SAC/AC

This message is set to be only visible to PC, SAC, AC, and the authors.

I want to disclose that I find reviewer A's evaluation of this paper quite unreasonable. This reviewer writes:

- Certain methods, such as method type, demonstrate limited effectiveness in setting, which may restrict their practical deployment.
- The paper points out that many of the an important task methods tested are essentially extensions of existing models adapted for another important task, such as a famous method.

It doesn't make much sense to cite the low performance of certain featured methods as weaknesses of a dataset-proposing/benchmark paper. It is not the authors' problem if an established method underperforms. Instead, the point of benchmarking is precisely to show when a method would fail. Many benchmark works have done this — some examples — and it is beyond comprehensible why this is considered a weakness.

Another criticism from reviewer A is:

- The tasks within the benchmark may not capture all possible real-world application scenarios, possibly overlooking specific needs within certain domains.

This, in my opinion, is a boilerplate concern that can be said for literally *any* dataset. While I do agree that the proposed dataset does not capture some important task scenarios — some examples — criticizing it for “not capturing all possible real-world applications” crosses the line and feels borderline hostile. This is akin to criticizing a method paper for not evaluating on every possible dataset.

I recommend the AC to either disregard A's review or consider encouraging the reviewer to revisit the evaluation.

851
852 This paper was ultimately accepted. This anecdote shows that without inner-reviewer analysis,
853 simple rule-based policies such as “inactive → desk rejection” fail to capture cases of severely low-
854 quality reviews. Finer-grained measures must be practiced to ensure that positive impact scales
855 broadly, rather than being limited to a few desk-rejected papers.
856
857858 C.2 CASE 2: REVIEWERS ASKING FOR PARTICULAR EXPERIMENTS AFTER THE REBUTTAL
859 DEADLINE.
860
861862 In a top ML conference where the exchange between authors and reviewers are limited to certain
863 time window, we had a split decision situation where the non/late-responding reviewers are not
supportive of the submission. As the AC is calling for consensus, we jumped in and asked:

864
865

Internal reviewer discussion

866
867
868
869
870
871
872

I skimmed over the two negative reviews of this work and found merits in many of the reviewer-raised points. However, I also find the authors' rebuttal to be proper in many regards — especially when the raised concern demands a clarification-like answer.

It looks like the two negative reviewers have yet to address the authors' rebuttal in a meaningful way (only acks are issued, cmiiw). So, to reach a consensus, I believe it would be helpful if the two reviewers could elaborate a bit on their leftover concerns. I am happy to set aside some time to discuss such leftover issues from my perspective.

873
874
875
876
877

Essentially, the two negative reviewers believed that certain experiments were missing — one of which can be seen as a combination of two existing methods, and another as a specific investigatory study of the author-proposed method. While we find such suggestions to have merit, we believe they were not raised appropriately from a procedural standpoint:

878
879

Internal reviewer discussion

880
881
882
883
884

I appreciate A's detailed response and updated review. I believe **A (as well as B)'s main concerns regarding ProposedMethod vs. PriorWork1 + PriorWork2 are legitimate and sound.** That being said, I am always the kind of reviewer who is "more in the authors' shoes"—for lack of better words—and I would like to present two alternative arguments regarding this concern.

885
886
887
888
889
890
891
892

First, I believe experiment-comparison requests that touch on *combinations of existing works* should be cautiously brought up. Many methods can be combined, but their combinations typically require a number of discretionary design decisions, and it is often unlikely for authors to feature the exact combination a reviewer has in mind. In this case, ProposedMethod proposes a paradigm of reducted, where the scope of eligible combinations is wide. Thus, in my opinion, **if reviewers are specifically interested in the comparison ProposedMethod vs. PriorWork1 + PriorWork2, such a request should be made explicitly before the rebuttal deadline, rather than mentioned in hindsight when the authors have no channel to address it.**

893
894
895
896
897
898
899
900
901
902
903
904
905

From the look of it, the ProposedMethod authors submitted their initial rebuttal on an early date, which is [redacted] days after the review post. However, only A engaged substantively on a late date. I must note that this year's BigConferenceName requires only two rounds of exchange, yet only 2/4 reviewers provided those to the authors, with all engaged reviewers leaning positive. **For such reasons, while I am also interested in this comparative result and agree with A's analysis, I do not believe we can use it against the authors (at least not as a singular veto reason), as the request was not properly raised from a procedural standpoint.** Imagine we were submitting a paper where reviewers were largely non-responsive, and the paper was then rejected for missing an experiment that was never explicitly asked for—it would be hard not to feel that is unfair. While I understand we all have different priorities and may have limited bandwidth for various reasons, we need to properly compensate authors when such situations occur.

906
907
908
909

(I know it is uncommon for a reviewer to argue on behalf of authors, but I always do so when it is warranted. The AC is welcome to confirm that this advocacy is from me for good reasons and not the result of any collusion.)

910
911
912
913
914

Later, we and the two reviewers exchanged more than five comments in total, which helped the AC reach a favorable decision. This anecdote shows that many reviewers are willing to engage in internal discussions, and such exchanges are profoundly helpful in deciding borderline papers — they simply never had the "push" to initiate such discussion voluntarily. We argue that our credit-based incentive system could help elicit more of these productive discussions.

915
916
917

Another takeaway from this exchange is the importance of having enforceable safeguards (e.g., reply deadlines), as otherwise authors may have no channel to meaningfully rebut at all. Attaching rewards and penalties to such actions is also crucial to ensure procedural fairness. While we respect and appreciate 'A' and 'B' for engaging our discussion, the fact that they were ok with not replying

918 / late replying authors is because there is virtually no penalty to them (as their "misconduct" would
 919 be too minute for desk rejection, yet the conference organizer has no other finer-grained tools) —
 920 something our credit system would help.
 921

922 C.3 CASE 3: AUTHORS PRESENT UNSUPPORTIVE RESULTS WHILE CLAIMING OTHERWISE. 923

924 In this case, we found that the authors were presenting unsupportive results to one reviewer while
 925 verbally claiming the opposite — so we stepped in:
 926

927 Internal reviewer discussion

928 As another reviewer, I find the reading on ProblemX tricky. The goal of Task is to have
 929 the [redacted].

930 In ProblemX, when the component is trained only on dataset1, the dataset2 ac-
 931 curacy drops quite significantly — a disadvantaged result. Once the component is trained
 932 on both dataset1 and dataset2:

- 933 • The dataset1 performance improves a small number over the a
 934 baseline, but this also comes at the cost of generating many more tokens than
 935 the a baseline.
- 936 • [redacted as it is too technical]
- 937 • In the end, the system shows only a small number accuracy gain on the task it
 938 was specifically trained on, while doing something at a higher cost.
 939 I think the added experiment makes the work more negative (though I appreciate
 940 the transparency), as it feels like the proposed component is very task-specific
 941 and does not generalize well.

942 (This message is set to be not visible to authors so that we can have a discussion with
 943 supposedly no bias, but I am happy to adjust the visibility should you want a more public
 944 discussion.)

945
 946 The reviewer exchanged views with us and we reached an agreement. The paper was ultimately
 947 rejected, with the AC explicitly citing our reasoning. Specifically, this other reviewer noted:
 948

949 Internal reviewer discussion

950 (BTW, you're among the very few reviewers who respond to other review comments. I truly
 951 admire this level of dedication.)

952
 953 This suggests that internal reviewer discussion is indeed rare, even though we lack direct statistical
 954 evidence.
 955

956 We hope these three anecdotal case studies illustrate how even a small component of our proposed
 957 framework — encouraging more internal reviewer discussion — can meaningfully facilitate paper
 958 evaluation. While we acknowledge our bias, we believe that in all three cases, the key factor behind
 959 each paper's acceptance or rejection was our initiative in starting those discussions. This suggests
 960 that most reviewers *want* to contribute and most ACs will take such discussions seriously; they
 961 simply need a small push to take the initiative — a push that our point-based incentive could very
 962 well provide.
 963
 964
 965
 966
 967
 968
 969
 970
 971