

Reading ability conflates SES creativity gaps

Guilherme Lichand^a, Leticia Lopes^b, and Sachin Allums^c

This manuscript was compiled on January 29, 2026

Who is more creative: high- or low-socioeconomic status (SES) individuals? This question is the focus of intense debates within both the scientific community and society at large, since creativity has been linked to innovation, productivity, and wealth generation. This question is, however, hard to answer; in particular, because creativity is typically assessed through standardized tests that build on reading and writing proficiency, which might conflate the relationship between SES and creativity. To overcome this challenge, we combine high-quality data on reading ability and experimental variation in reading requirements embedded in creativity assessments, in a series of studies with 6-12th graders in Brazilian schools. We first document that established creativity measures exhibit sizable SES gaps, but that these gaps become much smaller and no longer significant once we parse out students' reading ability (Study 1). Next, in two randomized control trials, we have students complete divergent thinking tasks while experimentally varying reading requirements: in the control group, students had to read question prompts, just like in standard assessments (e.g., PISA); in the treatment group, enumerators read the prompts on their behalf. In both conditions, enumerators wrote down students' answers. We find that high-SES students outperform low-SES ones *only* when they read the prompt, but not when performance does not depend on reading ability (Studies 2 and 3).

Creativity | SES gaps | Reading ability | Test unfairness

1. Introduction

Are high-SES individuals more creative than low-SES ones? There are compelling reasons to think this might be the case. Numerous studies find positive correlations between SES and creative performance (1–3), a pattern that holds even across cultures (4). SES tends to positively correlate with executive functions, which in turn, could have downstream effects on creative potential (5). SES could further boost subjects' confidence in their creative abilities, or their creative self-efficacy, which could then positively impact their creative skills (6). Creativity is often linked to innovation, entrepreneurship, productivity, and ultimately, wealth creation; thus, the most creative individuals may be those who accumulate wealth. Moreover, wealth might provide cognitive freedom, freeing people from immediate survival concerns and allowing them to think more imaginatively and beyond conventional constraints (3, 6, 7). In contrast, there are plausible mechanisms that might lead the exact opposite to hold. More creative individuals may feel unfulfilled by the repetitive tasks typical in many educational and workplace settings – with potential detrimental effects to their productivity, learning, and employability, particularly within societies that undervalue creative skills. Moreover, limited resources can push individuals to think creatively out of necessity, developing resourceful ways to meet basic needs (3, 8). Determining which effect dominates (if any) is an empirical question, widely debated in both academic research and public discourse.

The ascribed relevance of creative thinking, especially in a context of increasing automation since the advent of generative AI (9), has motivated recent efforts to measure creativity globally, such as the extensive Program for International Student Assessment (PISA) framework, which assesses 15-year-olds in over 60 countries. The 2022 exam specifically evaluated students' abilities to generate diverse and original

Significance Statement

We document that, while established creativity measures exhibit sizable socioeconomic status (SES) gaps, these gaps become much smaller and no longer significant once we parse out students' reading ability. Leveraging two randomized control trials, we find that high-SES students outperform low-SES ones in divergent thinking tasks *only* when assessments embed reading requirements.

Author affiliations: ^aStanford Graduate School of Education, Stanford, CA, United States. Corresponding author: glichand@stanford.edu.; ^bSão Paulo School of Economics (FGV-EESP), São Paulo, Brazil; ^cStanford University, Stanford, CA, United States

125 ideas, as well as to assess and refine ideas across various
126 contexts, using open-ended tasks focused on communication
127 and problem-solving. The OECD report on PISA creativity
128 scores documents a sizable socioeconomic status (SES) gap,
129 both across and within countries. Concretely, the gap between
130 the high-income OECD countries and Brazil (the setting of our
131 study) was 10 points in the PISA scale – 83% of a standard
132 deviation of OECD students' creative thinking performance.
133 This difference is comparable to the SES gap within Brazil,
134 where low-SES students scored 11 points lower than their high-
135 SES counterparts – a nearly 50% difference in average scores.
136 According to the report, a 1-unit increase in the SES index
137 within Brazilian 15-year-olds assessed by PISA corresponds to
138 a 4-point improvement in creative thinking performance (10).

139 At face value, these correlations support the theory that
140 creativity is positively correlated with wealth. However,
141 performing well on the PISA creativity assessment – a standard
142 assessment, based on open-ended questions – *also requires*
143 strong reading and writing skills. Since low-SES students
144 often have lower proficiency in these areas (3), the observed
145 SES gap in creativity scores may, at least in part, reflect
146 underlying disparities in reading and writing ability.

147 We report the results of three studies that address this
148 challenge, conducted with different samples of 6-12th graders
149 in Brazilian schools. Study 1 first documents SES gaps in
150 standard creativity measures and investigates the extent to
151 which reading ability mediates these gaps. Study 2 then
152 explores whether experimentally removing reading require-
153 ments from the creativity assessments used in Study 1 closes
154 these SES gaps. Finally, Study 3 replicates that experimental
155 manipulation in the context of items from the PISA 2022
156 Creative Thinking assessment.

157 This paper hopes to establish whether SES gaps in creativity
158 measures can be attributed to reading requirements rather
159 than to meaningful differences in creative potential. This work
160 adds to an ongoing debate about SES gaps in creativity, with
161 implications for the assessment of higher-order skills in general,
162 and creativity in particular.

164 2. Results

165 **Study 1: SES creativity gaps in standard creativity measures**
166 **conflate reading ability.** In line with evidence from interna-
167 tional assessments like PISA and with previous research
168 findings, Figure 1 documents a sizable and statistically
169 significant SES creativity gap based on the summary measure
170 of conventional creativity measures (the AUT and the DAT).
171 Concretely, high-SES students score 0.154 s.d. higher than
172 low-SES students (p-value of difference = 0.023).

173 Figure 2 documents that, however, this gap at least partly
174 conflates differences in reading ability across low- and high-SES
175 participants. Leveraging high-quality reading ability data from
176 ROAR-SRE, we document that while the SES creativity gap
177 among proficient students is of similar magnitude (0.150 s.d.),
178 albeit imprecisely estimated (p-value of difference = 0.076),
179 such gap is 1/3 lower among those who do not read at grade-5
180 level (0.098 s.d.) and no longer statistically significant (p-value
181 of difference = 0.324). Supplementary Materials document
182 that results hold for the different components of the summary
183 measure.

184 It is also striking that, even among students who read
185 proficiently, the SES creativity gap is 2/3 lower than the

187 creativity gap across reading ability levels in the whole sample
188 (0.470 s.d., p-value of difference < 0.001).

189 Together, the patterns we document are consistent with
190 the claim that reading ability conflates SES creativity gaps in
191 standard creativity measures.

192 Nonetheless, since reading ability is not randomly assigned
193 – but, rather, correlated with several student characteristics,
194 including household wealth –, it could be that our findings
195 indicate that our reading ability measure is merely a better
196 proxy for SES than the one we rely on (number of full
197 bathrooms). For this reason, Studies 2 and 3 introduce
198 experimental variation in reading requirements embedded in
199 creativity assessments.

200 **Study 2: High-SES students outperform low-SES ones in** 201 **standard creativity measures only when they read the prompts.**

202 Figure 3 documents that high-SES students outperform low-
203 SES ones in standard creativity measures (the AUT and the
204 DAT) only due to reading requirements.

205 Among participants who had to read the question prompts,
206 we estimate a SES creativity gap similar to the unconditional
207 one documented in Study 1 (0.134 s.d., p-value of the difference
208 = 0.023). Conversely, among those for whom enumerators
209 read the prompts on their behalf, the SES gap is less than 1/3
210 (0.068 s.d.) and no longer statistically significant (p-value of
211 the difference = 0.249). Supplementary Materials document
212 that results hold for the different components of the summary
213 measure.

214 The experimental findings provide causal evidence that
215 SES gaps in conventional creativity measures indeed conflate
216 differences in reading ability.

217 Next, Study 3 investigates whether the same holds for
218 less conventional creativity measures included in international
219 assessments like PISA.

220 **Study 3: High-SES students outperform low-SES ones** 221 **in PISA-like creativity measures only when they read the** 222 **prompts.**

223 Figure 4 documents that, similar to Study 2, reading
224 requirements mask differences between high- and low-SES
225 students in PISA-like creativity measures (diverse and creative
226 idea generation tasks).

227 Among participants who had to read the question prompts,
228 we estimate virtually no SES creativity gap (a 0.004 s.d.
229 difference in favor of high-SES students, p-value = 0.903).
230 In turn, among those for whom enumerators read the prompts
231 on their behalf, a sizeable SES gap emerges *in favor of low-SES*
232 *students* (0.085 s.d., p-value = 0.011).

233 Supplementary Materials replicate the analyses for the
234 different components of the summary measure. Consistent
235 with the main text, SES differences in the share of students
236 who receive full credit in the first item (based on the PISA's
237 scoring guidelines) are reversed in the absence of reading
238 requirements (p-value of difference-in-difference = 0.043).

239 Supplementary Materials further show that, when using
240 GenAI instead of hand-coded scores, the patterns of SES
241 gaps match precisely those documented in Study 2: high-SES
242 students *only* outperform low-SES ones in the presence of
243 reading requirements.

244 3. Discussion

245 Together, our findings provide compelling evidence that
246 reading ability conflates creativity measures in standard assess-
247
248

249 ments – including PISA 2022 –, in particular when it comes
250 to SES gaps. The OECD report acknowledges that “[a]fter
251 accounting for students’ mathematics and reading performance,
252 differences in the performance of advantaged and disadvan-
253 taged students are much smaller in all countries/economies –
254 and even become statistically non-significant in 14 [of them]”
255 (11, p. 113). While it posits that “[S]ocio-economic disparities
256 in creative thinking performance therefore rather reflect a
257 range of economic and cultural factors, experiences and
258 mechanisms known to affect student achievement overall” (11,
259 p. 113), it does not mention test unfairness as a potential driver
260 of SES differences. Instead, the report concludes that, since
261 the correlation between reading and creativity scores across
262 OECD countries is 0.66 – less than that between reading and
263 other disciplines –, “the creative thinking assessment measures
264 a different subset of skills with respect to those measured in
265 the mathematics, reading and science assessments” (11, p. 83).
266 This study challenges this conclusion; concretely, it suggests
267 that SES differences in reading comprehension can account
268 for the *totality* of SES differences in average creativity scores
269 across income groups within countries, and across countries of
270 different income levels.

271 More broadly, while previous research documents a positive
272 correlation between SES gaps in literacy skills and those in
273 other subjects, like math (12), it is hard to infer unfairness
274 solely based on these associations. As (13) points out,
275 “[a] major limitation of item bias statistics or indices is
276 that measures of relative difficulty do not provide proof of
277 unfairness. Only if an item is relatively more difficult for one
278 group (statistically biased) *and* the source of this difficulty is
279 irrelevant to the test construct is an item said to be unfair”
280 (13, p. 234). Our paper innovates by experimentally varying
281 whether test performance hinges on reading proficiency –
282 a source of difficulty irrelevant to the divergent thinking
283 constructs that the PISA aims to capture. Having enumerators
284 write down responses for participants further allows us to
285 isolate creativity from writing ability. All in all, we provide
286 causal evidence that the control version of the test is *unfair*.

287 An active literature (e.g., 14, 15) discusses how creativity
288 might reflect a range of underlying factors, from general
289 intelligence to other latent variables. If both reading ability
290 and creativity reflect cognition more broadly (construct
291 redundancy), one could ask whether it is surprising at all
292 that when we control for one (as in Study 1) we partially
293 eliminate variation in the other. This is not, however, what
294 we do in Studies 2 and 3, which experimentally isolate the
295 role of reading requirements – not that of reading ability or
296 its underlying factors – in expressing creative ability.

297 Our findings are naturally limited by our study design.
298 Even if Studies 1 and 2 conducted the same tests, differences in
299 student population, assessment type, and time limits constrain
300 our ability to directly compare SES gaps across them. In
301 particular, Study 1 was self-administered: participants had
302 2 minutes to complete the AUT and 4 minutes to complete
303 the DAT (computer-based) as they had to write down their
304 responses. In contrast, Studies 2 and 3 were mediated
305 by surveyors: across experimental conditions, participants
306 had only 30 seconds to complete the AUT and 45 seconds
307 to complete the DAT (tablet-based) because enumerators
308 wrote responses on their behalf. Despite these differences in
309 administration and sampling (a non-random sample of Rio
310

311 middle-schoolers in Study 1, and a representative sample of
312 Brazilian K-12 students in Studies 2 and 3), we observe SES
313 creativity gaps of very similar magnitudes in the AUT and
314 DAT in both cases.

315 We carefully designed the three studies to ensure our
316 findings are not country-specific. Studies 1 and 2 featured
317 the AUT, the most widely used divergent thinking task in
318 the creativity literature (16), and the DAT, a task with an
319 objective scoring methodology, unlike many other divergent
320 thinking tasks (16, 17). We chose these measures in particular
321 because they have been traditionally used as evidence that
322 high-SES subjects might be more creative than low-SES ones
323 (3). The fact that we find similar patterns across the lab-in-
324 the-field AUT and DAT tests in Studies 1 and 2 and the PISA
325 items in Study 3 (which presumably align more closely to real-
326 world creative problem solving) lends some confidence that our
327 findings are relevant above and beyond the artificial setting
328 of our experiments. While we use fluency as the primary
329 measure of creativity on the AUT, Supplementary Materials
330 show that our conclusions extend to other ways of scoring the
331 AUT, such as originality, which might more closely measure
332 true creative potential instead of merely conflating generating
333 ideas with creative ability (18). The items used in Study
334 3 were extracted from PISA itself, and the SES creativity
335 gaps in our control group correspond almost exactly to those
336 documented in the OECD report. Brazil was the setting of
337 this study thanks to the unique opportunity to conduct large-
338 scale survey experiments in more than 200 schools throughout
339 the country – a research-practice partnership that is hard to
340 replicate anywhere else.

341 Despite these efforts, replicating these experiments in other
342 settings and with additional creativity measures to ensure that
343 these results generalize further remains a promising avenue
344 for future work. Comparative education research directly
345 contrasting the performance of high-income OECD countries
346 and that of lower-income countries like Brazil, featuring the
347 experimental manipulation of Studies 2 and 3, could provide
348 convincing evidence to challenge the findings of the OECD
349 report. It also remains unclear if these findings would extend
350 to other aspects of verbal creativity (i.e., convergent thinking),
351 other assessments for divergent thinking (i.e., form-first vs
352 function-first tasks), or even figural creativity (completing
353 abstract figures in a novel way; e.g., 19). As a leading example,
354 while reading instructions aloud could help in figural creativity
355 tasks, since participants might understand more clearly that
356 they should not simply complete an obvious image, we cannot
357 say for certain whether that would be the case. Future work
358 could focus on using additional measures to examine whether
359 reading ability also conflates SES creativity gaps in those
360 cases.

361 Overall, our findings suggest that assessments of higher-
362 order skills, such as creative thinking, should better account for
363 limitations in foundational skills (often disproportionately con-
364 centrated among low-income students). Specifically, assessors
365 might consider administering oral exams or even non-verbal
366 measures of creativity or other higher-order skills in order to
367 accurately capture underlying skills. More broadly, in line
368 with factor-analytic frameworks (e.g., 14, 15, 20), our findings
369 suggest that reading requirements embedded in conventional
370 assessments might conflate a host of cognitive abilities, from
371 executive functions to mathematical reasoning. Ultimately,
372

373 unfairness in test design not only constrains the extent to which
374 the assessment accurately reflects levels and differences in the
375 skills or psychological constructs it intends to capture, but also,
376 might end up reinforcing deficit narratives that mischaracterize
377 the potential of disadvantaged students.
378

379 Methods

380 Study 1 was approved by the Stanford University Institutional
381 Review Board (protocol #77799) on May 29, 2025. Studies
382 2 and 3 were approved by the Brazilian Institutional Review
383 Board (CONEP; protocol #7.168.931) on October 18, 2024,
384 and by the Stanford University Institutional Review Board
385 (protocol #77912) on January 17, 2025. Waiver of parental
386 consent obtained from all school principals, as the studies
387 involved no more than minimal risk to study participants,
388 whom still had to assent and could drop out of the study
389 at any point and without penalty. Supplementary Materials
390 contain additional descriptive statistics information on the
391 study samples and robustness checks for our findings.
392

393 **Study 1: Documenting SES creativity gaps and the extent to**
394 **which reading ability mediates them.** In May-Jul/2025, the
395 study enrolled 523 students across 8 middle schools (grades
396 6-9) in the municipality of Rio de Janeiro, Brazil. Recruited
397 students participated in computer-based reading and creativity
398 assessments, covering conventional divergent thinking tasks.
399 Divergent thinking is a core component of creative thinking:
400 the ability to generate a variety of creative possibilities and
401 novel associations (21, 22).
402

403 The analyses proxy for SES by comparing students living in
404 households with at most one full bathroom (*low-SES*) to those
405 with two or more (*high-SES*). This is considered a reliable
406 proxy for SES in contexts where direct measures of income
407 or wealth are unavailable or impractical (23). Since children
408 were the ones filling out the survey with no caregivers present,
409 we decided to use full bathrooms as a proxy for SES. Based
410 on this proxy, 366 study participants were low-SES and 157
411 were high-SES – a distribution very similar to the Brazilian
412 national figures, based on the 2022 Census (24).

413 We measure students' reading ability through the Rapid
414 Online Assessment of Reading (ROAR; 25–29), a digital
415 assessment that measures reading fluency and comprehension.
416 Concretely, we assessed participants' reading fluency and
417 comprehension using the ROAR Sentence Reading Efficiency
418 (ROAR-SRE) task, previously validated for Brazil (29). In
419 the task, participants rate sentences as true or false as
420 quickly and accurately as they can, within 180 seconds.
421 Sentences were carefully validated for Brazilian Portuguese
422 and screened for contextual suitability, requiring minimal
423 background knowledge and using simple vocabulary and
424 syntactic structures (29). Sentences are presented in random
425 order. We score the task as the number of correct answers
426 minus the number of incorrect answers over that time span,
427 in line with (29).

428 Based on (29), which associates ROAR-SRE scores to
429 proficiency labels in Brazil, 43.02% (N = 225) of 6-9th graders
430 in our sample do not read at grade-5 level (43.44% among
431 low-SES students, and 42.04% among high-SES ones).

432 We conducted five trials of the Alternative Uses Task (AUT)
433 and one trial of the Divergent Association Task (DAT) with
434 all study participants in the context of a companion paper

(30), which randomly assigned a generative AI tool to assist
students with the creativity tasks. In this study, we restrict
attention to the AUT unassisted trials (the first and the last
ones), and parse out treatment indicators across all analyses.

In the AUT, participants have 2 minutes to come up with
as many creative and unusual uses for an everyday object.
Participants completed the AUT for the five following items:
tire, pants, shoe, table, and bottle. These objects were chosen
because their uses have the highest correlation of judgements
between either an automated or human evaluator for originality
(31). The order of the objects was randomly assigned across
participants.

The AUT is typically scored on different dimensions:
fluency (the number of distinct uses), *flexibility* (the number
of categories to which alternative uses correspond), and
originality (responses' statistical infrequency) – each deemed
as a key indicator of creative potential (21, 22, 32, 33). In
the main text, we restrict attention to fluency, the most
common and transparent dimension used in the literature;
Supplementary Materials replicate all analyses using originality
instead of fluency.

In the DAT, participants have 4 minutes to come up with
10 nouns as semantically different as possible from one another
(17). For example, 'dog' and 'cat' are very close in meaning,
and thus, including both of them would result in a low DAT
score. By contrast, 'dog' and 'guitar' would score much higher
since they have very different meanings and uses. In the DAT,
words must be single words, non-proper nouns, and not be
considered technical vocabulary. Only the first seven valid
words are scored, following the scoring algorithm in (17). If a
participant did not provide seven valid words, they received a
score of zero.

Participants could refuse answering one or more trials;
response rates for each trial were statistically identical across
SES (92.1% for low-SES and 91.7% for high-SES, p-value of
difference = 0.848).

We compute p-values from Wald tests of differences in
average test scores between high- and low-SES participants, at
first, unconditional on reading ability, and next, conditioning
on whether participants could read at grade-5 level (*high*
literacy) or not (*low literacy*). To deal with family-wise error
rates in hypotheses testing from multiple comparisons, we focus
on a divergent thinking *summary measure*, stacking the dataset
with participants' test scores in each trial (2 observations
for AUT and 1 for DAT) and clustering the analyses at the
student level. We also control for age, gender, and treatment
assignment in the companion paper. Supplementary Materials
showcase results without controls, and for each component of
the summary measure.

Study 2: Experimentally parsing out reading ability in es-
tablished creativity measures. In May-Jun/2025, the study
enrolled 1,100 participants across 140 high schools (grades
10-12) in 27 Brazilian States. The sample was drawn to
represent all K-12 schools in the country. Recruited students
participated in tablet-based creativity assessments, adapted
from those conducted in Study 1.

We again proxy for SES by comparing students living in
households with at most one full bathroom (*low-SES*) to those
with two or more (*high-SES*). Based on this proxy, 681 study
participants are low-SES, and 419, high-SES.

497 We adapted the same AUT and DAT tasks from Study 1 to
498 be conducted on tablets in a representative sample of schools
499 all over the country. Due to time limitations, we ran only two
500 trials of the AUT (all subjects were prompted with a tire and
501 a bottle, in that order), 30 seconds each, and one trial of the
502 DAT, for 45 seconds.

503 All subjects received a seed word for the DAT ('ball'),
504 which was not considered in their final score. Across all condi-
505 tions, participants said their answers aloud, and enumerators
506 recorded their responses.

507 We independently randomized whether each question
508 prompt was read by the student (control group) or by the
509 enumerator (treatment group). As such, in the control group,
510 our measure of creativity may still reflect students' reading
511 ability, while in the treatment group, it isolates creativity from
512 reading skills. Across both treatment and control conditions,
513 students responded verbally, and enumerators typed their
514 responses; as a result, the creativity measure is independent of
515 students' writing skills in both groups (see SI for all details).

516 In the AUT, the control group includes 559 participants,
517 and treatment group, 541; in the DAT, the control group
518 includes 546 participants, and the treatment group, 554.
519 Participants' characteristics were balanced across control and
520 treatment conditions (p-value of equality across all variables
521 and conditions = 0.392).

522 Like Study 1, participants could refuse answering one or
523 more trials, and response rates for each trial were statistically
524 identical across SES (95.8% for low-SES and 96.5% for high-
525 SES, p-value of difference = 0.442).

526 We compute p-values from Wald tests of differences in
527 average test scores between high- and low-SES participants
528 conditioning on whether participants read the prompt them-
529 selves or whether enumerators read them on their behalf. As in
530 Study 1, we focus on a divergent thinking summary measure,
531 stacking the dataset with participants' test scores in each
532 trial (2 observations for AUT fluency and 1 for DAT) and
533 clustering the analyses at the student level. We also control
534 for age, gender, and location (= 1 if urban, 0 otherwise).
535 Supplementary Materials showcase additional results: without
536 controls, for each component of the summary measure, and
537 using AUT originality instead of fluency.

538 **Study 3: Experimentally parsing out reading ability in PISA**
539 **measures.** In Aug-Sep/2024, the study enrolled 515 partic-
540 ipants across 43 high schools in 23 Brazilian States. The
541 sample was drawn to represent all K-12 schools in the country.
542 Recruited students participated in tablet-based creativity
543 assessments, adapted from the PISA 2022 Creative Thinking
544 Assessment.

545 We once again proxy for SES by comparing students living
546 in households with one full bathroom (*low-SES*) to those
547 with two or more (*high-SES*). Based on this proxy, 266 study
548 participants are low-SES, and 249, high-SES.

549 In the experiment, we adapted two divergent thinking tasks
550 from the PISA 2022 Creative Thinking Assessment framework.
551 Each item assesses distinct ideation processes. The first item
552 targets generating *diverse ideas*, emphasizing students' ability
553 to think flexibly. The second item relates to generating *creative*
554 *ideas*, focusing on the novelty and usefulness of an idea.

555 Specifically, the items were sourced from the Written
556 Expression domain of the PISA 2022 Creative Thinking
557 Assessment. In the first item, "Generate Diverse Ideas",
558

559 high-school students are tasked with creating three distinct
560 titles based on a surreal image displayed on a tablet, which
561 features a large-scale book alongside various elements of the
562 countryside, such as a field, tree, and wooden bench. They
563 are given two minutes to complete this task. The second
564 item, "Generate Creative Ideas", presents a dialogue between
565 the sun and the earth in a comic strip format. This exercise
566 consists of six dialogue boxes, whereby students must craft
567 sentences for each character and communicate their responses
568 to the enumerator within three minutes. These ideation
569 processes embody cognitive skills essential for creative thinking
570 in educational settings (17, 32, 34).

571 Responses were evaluated based on PISA's criteria for
572 measuring creativity levels. Each item is associated with a
573 specific difficulty level that determines the maximum score
574 attainable in that item. A response can be rated as full credit,
575 partial credit, or no credit. We resorted to two approaches
576 for scoring students' responses. First, generative AI (GenAI),
577 providing ChatGPT with item-specific prompts that include
578 the evaluation criteria, classification parameters, and examples
579 of responses associated with each credit score derived from the
580 PISA materials. This approach leverages GenAI's capacity to
581 analyze and interpret data, allowing us to evaluate student
582 responses systematically (see SI for all details). Second, we
583 had research assistants hand-code responses following the same
584 prompts (see SI for all details). In both cases, our analysis
585 focuses on the percentage of responses rated as *full credit*,
586 based on the PISA scoring guidelines.

587 We once again independently randomized whether each
588 question prompt was read by the student (control group) or by
589 the enumerator (treatment group). Consistent with Study 2,
590 students responded verbally across both treatment and control
591 conditions, and enumerators typed their responses (see SI for
592 all details).

593 In the first item, the control group includes 264 participants,
594 and treatment group, 251; in the second item, the control
595 group includes 243 participants, and the treatment group, 272.
596 Participants' characteristics were balanced across control and
597 treatment conditions (p-value of equality across all variables
598 and conditions = 0.355).

599 As before, participants could refuse answering one or more
600 trials; response rates for each trial were statistically identical
601 across SES (80.8% for low-SES and 81.3% for high-SES, p-
602 value of difference = 0.857).

603 We compute p-values from Wald tests of differences in
604 average test scores between high- and low-SES participants
605 conditioning on whether participants read the prompt them-
606 selves or whether enumerators read them on their behalf. As
607 in Studies 1 and 2, we focus on a divergent thinking summary
608 measure, stacking the dataset with participants' test scores
609 in each trial (one for each task) and clustering the analyses
610 at the student level. We also control for age, gender, and
611 location, in line with Study 2. In the main text, we focus on
612 hand-coded scores. Supplementary Materials present results
613 without controls, for each component of the summary measure,
614 and using GenAI scores instead of hand-coded ones.

615 Data Availability

616 All data are deposited at Open Science Framework (<https://osf.io/d7hm8/>).

Code Availability

A complete replication package is deposited at Open Science Framework (<https://osf.io/d7hm8/>).

Acknowledgements

Data collection for Study 1 conducted in partnership with the Municipal Secretariat of Education of Rio de Janeiro. Data collection for Studies 2 and 3 conducted by Equidade.info, generously funded by the Stanford Lemann Center. We acknowledge helpful discussions with Ben Domingue, Marilyn Oppezzo and Dan Schwartz, and excellent research assistance by João Bigon, Lucas Garcia, Paulina Huang, Rodrigo Megale, Bianca Mesquita, Jônatas Ribeiro, Karine Roncete and Mavigson Silva. All remaining errors are ours.

Author Contributions

G.L. conceptualized and supervised the study and wrote the manuscript. L.L. and S.A. developed the code, analyzed the data, and wrote the manuscript. All authors contributed to reporting the findings.

Competing Interests

None of the authors have competing interests to disclose.

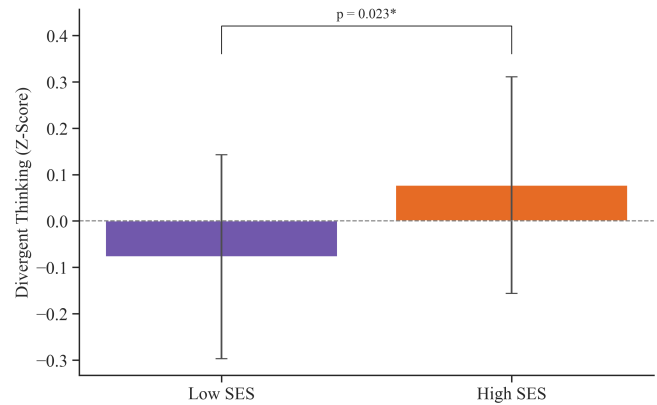
References

1. Mauricio Castillo-Vergara, Nicole Barrios Galleguillos, Laura Jofré Cuello, Alejandro Alvarez-Marín, and Christian Acuña-Opazo. Does socioeconomic status influence student creativity? *Thinking Skills and Creativity*, 29:142–152, 2018. ISSN 1871-1871. . URL <https://www.sciencedirect.com/science/article/pii/S1871187118300361>.
2. M.R. Sarsani. Socio-Economic Status and Performance on Creativity Tests. In *Encyclopedia of Creativity*, pages 360–363. Elsevier, 2011. ISBN 978-0-12-375038-9. . URL <https://linkinghub.elsevier.com/retrieve/pii/B9780123750389001795>.
3. Selcuk Acar, Harun Tadik, Recep Uysal, Danielle Myers, and Betül Inetas. Socio-economic status and creativity: A meta-analysis. *The Journal of Creative Behavior*, 57(1):138–172, 2023.
4. Zhitian Skylor Zhang, Linda Hoxha, Abdullah Aljughaiman, Aliriza Arënlju, Maria P. Gomez-Arizaga, Sule Gucyeter, Irina Ponomareva, Jiannong Shi, Paula Irueste, Silke Rogl, Miguelina Nunez, and Albert Ziegler. Social Environmental Factors and Personal Motivational Factors Associated with Creative Achievement: A Cross-Cultural Perspective. *The Journal of Creative Behavior*, 55(2):410–432, 2021. ISSN 2162-6057. . URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/jocb.463>.
5. Karina Hendrie Kupczynszyn, Vanessa Arán Filippetti, and Laura Oros. Socioeconomic status effects on children's creativity: The mediating role of executive functions. *Thinking Skills and Creativity*, 51:101437, 2024. ISSN 1871-1871. . URL <https://www.sciencedirect.com/science/article/pii/S1871187123002055>.
6. Zhongyao Lu, Yaqiong Ding, and Yanhong Nie. How does family socioeconomic status affect creativity? the role of creative self-efficacy and critical thinking disposition. *Current Psychology*, 43(6):5674–5681, 2024.
7. David Yun Dai, Xiaoyuan Tan, Deepti Marathe, Anna Valtcheva, Robert M Pruzek, and Jiliang Shen. Influences of social and educational environments on creativity during adolescence: Does sex matter? *Creativity Research Journal*, 24(2-3):191–199, 2012.
8. Sofie Dahlman, Per Bäckström, Gunilla Bohlin, and Örjan Frans. Cognitive abilities of street children: Low-SES Bolivian boys with and without experience of living in the street. *Child Neuropsychology*, 19(5):540–556, 2013.
9. Janet Ratner, Roger E Beaty, James C Kaufman, Todd Lubart, and Jacob Sherson. Creativity in the age of generative ai. *Nature Human Behaviour*, 7(11):1836–1838, 2023.
10. OECD. *Student performance in creative thinking*. OECD Publishing, Paris, 2024.
11. OECD. Pisa 2022 creative thinking, 2022. URL <https://www.oecd.org/en/topics/sub-issues/creative-thinking/pisa-2022-creative-thinking.html>. Accessed: 2024-10-10.
12. Anna Gomez, Elena Pecina, Sara Villanueva, and Tanya Huber. The undeniable relationship between reading comprehension and mathematics performance. *Issues in Educational Research*, 30(4):1329–54, 2020.
13. Gregory Camilli. Test fairness. In *Educational measurement (4th ed.)*, pages 221–256. American Council on Education/Praeger, 2006.
14. J. B. Carroll. *Human Cognitive Abilities: A Survey of Factor-Analytic Studies*. Cambridge University Press, 1993.
15. Jonathan Wai and Matt Brown. Developmental histories facilitating the emergence of creative scientific expertise: The role of developed cognitive talents, education, and social and cultural contexts. *Front. Psychol.*, 12, 2021. . URL <https://doi.org/10.3389/fpsyg.2021.716529>.
16. Janika Saretzki, Boris Forthmann, and Mathias Benedek. A systematic quantitative review of divergent thinking assessments. *Psychology of Aesthetics, Creativity, and the Arts*, 2024. .
17. Jay A Olson, Johnny Nahas, Denis Chmoulevitch, Simon J Cropper, and Margaret E Webb. Naming unrelated words predicts creativity. *Proceedings of the National Academy of Sciences*, 118(25):e2022340118, 2021.
18. Paul J. Silvia, Beate P. Winterstein, John T. Willse, Christopher M. Barona, Joshua T. Cram, Karl I. Hess, Jenna L. Martinez, and Crystal A. Richard. Assessing creativity with divergent thinking tasks: Exploring the reliability and validity of new subjective scoring methods. *Psychology of Aesthetics, Creativity, and the Arts*, 2(2):68–85, 2008. .
19. Andreas Fink, Thomas Reim, Mathias Benedek, and Roland H. Grabner. The Effects of a Verbal and a Figural Creativity Training on Different Facets of Creative Potential. *The Journal of Creative Behavior*, 54(3):676–685, 2020. . URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/jocb.402>.
20. Douglas Bors. The factor-analytic approach to intelligence is alive and well: A review of Carroll, J.B. (1993). Human cognitive abilities: A survey of factor-analytic studies. *Canadian Journal of Experimental Psychology*, 47:763–766, 1993.
21. J. P. Guilford. Creativity. *American Psychologist*, 5(9):444–454, 1950. . URL <https://doi.org/10.1037/h0063487>.
22. Mark A Runco and Selcuk Acar. Divergent thinking as an indicator of creative potential. *Creativity research journal*, 24(1):66–75, 2012.
23. Young J Juhn, Timothy J Beebe, Dawn M Finnie, Jeff Sloan, Philip H Wheeler, Barbara Yawn, and Arthur R Williams. Development and initial testing of a new socioeconomic status measure based on housing data. *Journal of Urban Health*, 88:933–944, 2011.
24. Igor Ferreira. Censo 2022: Rede de esgoto alcança 62,5% da população, mas desigualdades regionais e por cor e raça persistem | Agência de Notícias, 2024. URL <https://agenciadenoticias.ibge.gov.br/agencia-noticias/2012-agencia-de-noticias/noticias/39237-censo-2022-rede-de-esgoto-alcanca-62-5-da-populacao-mas-desigualdades-regionais-e-por-cor-e-ra>
25. J. D. Yeatman, K. A. Tang, P. M. Donnelly, M. Yablonski, M. Ramamurthy, I. I. Karipidis, S. Caffarra, M. E. Takada, K. Kanopka, M. Ben-Shachar, and B. W. Domingue. Rapid online assessment of reading ability. *Scientific Reports*, 11(1):6396, 2021. . URL <https://doi.org/10.1038/s41598-021-85907-x>.
26. J. D. Yeatman, J. E. Tran, A. Burkhardt, W. A. Ma, J. Mitchell, M. Yablonski, and A. Richie-Halford. Development and validation of a rapid online sentence reading efficiency assessment, December 1 2023. URL <https://doi.org/10.31219/osf.io/u3mjz>. Preprint.
27. J. D. Yeatman, C. Townley-Flores, K. Wentzlof, W. A. Ma, J. M. Siebert, M. Fuentes-Jimenez, A. Saavedra, T. S. Murray, K. Bhat, M. Ramamurthy, and ROAR Developer Consortium. Rapid online assessment of reading (roar): Technical manual. Technical report, Stanford

683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
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

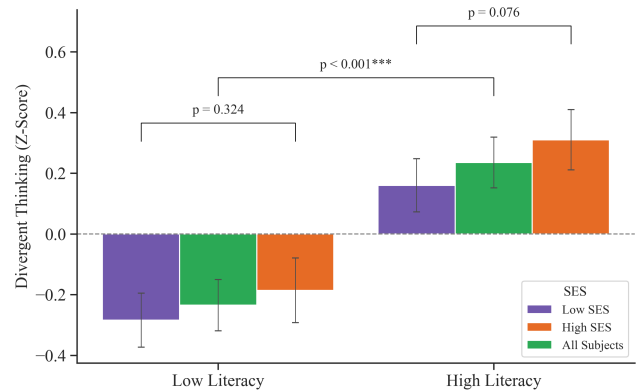
745 GSE, November 5 2024.
 746 28. W. A. Ma, A. Richie-Halford, A. Burkhardt, K. Kanopka, C. Chou, B. Domingue, and J. D.
 747 Yeatman. Roar-cat: Rapid online assessment of reading ability with computerized adaptive
 748 testing. *Working Paper*, 2023.
 749 29. Karine Roncete, Lucas Klotz, Wanjing Ma, Emily Arteaga, Luciana Alves, Rebecca Chrispin,
 750 Danielle Diniz, Jason Yeatman, and Guilherme Lichand. Development and validation of a
 751 rapid online sentence reading efficiency assessment, 2024. URL
 752 <https://www.researchsquare.com/article/rs-5516837/v1>. Preprint.
 753 30. Sachin Allums, Karine Roncete, and Guilherme Lichand. Generative AI can harm learning
 754 despite guardrails: Evidence from middle-school creativity. *mimeo*, 2025.
 755 31. Peter Organisciak, Selcuk Acar, Denis Dumas, and Kelly Berthiaume. Beyond semantic
 756 distance: Automated scoring of divergent thinking greatly improves with large language
 757 models. *Thinking Skills and Creativity*, 49:101356, September 2023. ISSN 1871-1871. . URL
 758 <https://www.sciencedirect.com/science/article/pii/S1871187123001256>.
 759 32. OECD. *Measuring creative thinking*. OECD Publishing, Paris, 2024.
 760 33. Joy Paul Guilford. The structure of intellect. *Psychological bulletin*, 53(4):267, 1956.
 761 34. Brian J Lucas and Loran F Nordgren. The creative cliff illusion. *Proceedings of the National
 762 Academy of Sciences*, 117(33):19830–19836, 2020.

807 **Figures**



825 **Fig. 1.** SES creativity gap in divergent thinking summary measure

826 **Notes:** Summary measure of two trials of the Alternative Uses Task (AUT) and one
 827 trial of the Divergent Association Task (DAT), stacking datasets (i.e., each observation
 828 is a student \times trial) while standardizing the outcome by the mean and standard
 829 deviation of the whole sample within each trial. Low-SES participants report to live in
 830 households with at most 1 full bathroom; high-SES ones, with 2 or more full bathrooms.
 831 N = 523 (366 low-SES and 157 high-SES participants). P-value of a Wald test of
 832 difference in means, clustering standard errors at the participant level.
 833



834 **Fig. 2.** SES creativity gaps, by reading ability

835 **Notes:** Summary measure of two trials of the Alternative Uses Task (AUT) and one
 836 trial of the Divergent Association Task (DAT), stacking datasets (i.e., each observation
 837 is a student \times trial) while standardizing the outcome by the mean and standard
 838 deviation of the whole sample. Low-SES participants report to live in households with
 839 at most 1 full bathroom; high-SES ones, with 2 or more full bathrooms. High (low)
 840 literacy participants are those who do (not) read at grade-5 level based on ROAR-SRE
 841 (see 29). N = 523 (366 low-SES and 157 high-SES participants). P-value of Wald
 842 tests of differences in means within reading ability level, clustering standard errors at
 843 the participant level.
 844

869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930

931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992

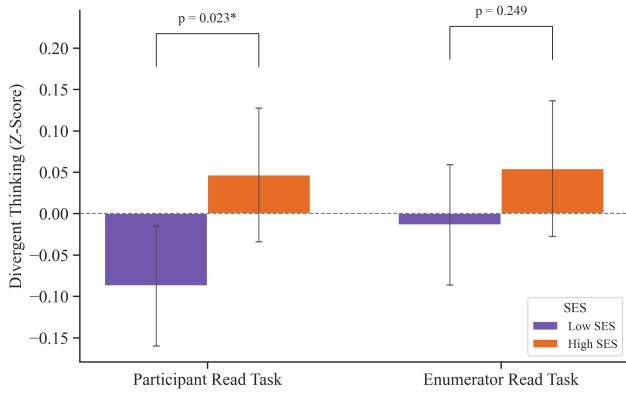


Fig. 3. SES creativity gaps in standard measures, by experimental condition

Notes: Summary measure of two trials of the Alternative Uses Task (AUT) and one trial of the Divergent Association Task (DAT), stacking datasets (i.e., each observation is a student \times trial) while standardizing the outcome by the mean and standard deviation of the control group (in which participants read the prompt) within each trial. Low-SES participants report to live in households with at most 1 full bathroom; high-SES ones, with 2 or more full bathrooms. $N = 1,100$ (681 low-SES and 419 high-SES participants). P-value of Wald tests of differences in means within each experimental condition, clustering standard errors at the participant level.

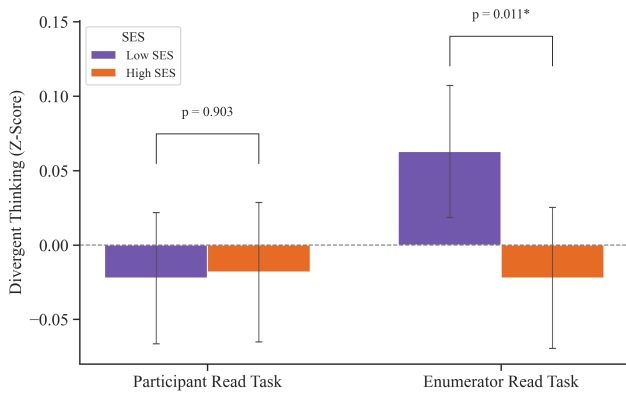


Fig. 4. SES creativity gaps in PISA-like items, by experimental condition

Notes: Summary measure of two items adapted from the PISA 2022 Creative Thinking Assessment, stacking datasets (i.e., each observation is a student \times trial) while standardizing the outcome by the mean and standard deviation of the control group (in which participants read the prompt) within each trial. Low-SES participants report to live in households with at most 1 full bathroom; high-SES ones, with 2 or more full bathrooms. $N = 515$ (266 low-SES and 249 high-SES participants). P-value of Wald tests of differences in means within each experimental condition, clustering standard errors at the participant level.