

---

# Potential of LLM-Generated Lifestyle Adjustment Recommendations Based on Multimodal Data

---

**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

1 The prevalence of burnout, depression, and stress-related disorders has increased  
2 markedly in contemporary societies, particularly in the context of flexible and remote  
3 working arrangements. These structural shifts impose novel demands on individuals  
4 to self-regulate health, well-being, and productivity—responsibilities that were  
5 previously supported by organizational structures in conventional workplaces.  
6 Traditional self-management strategies struggle to address the complexity of  
7 interacting behavioral, psychological, and physiological determinants. This paper  
8 explores the feasibility of employing large language models (LLMs) to generate  
9 lifestyle adjustment recommendations based on multimodal data that integrate  
10 subjective self-reports with objective sensor-derived measures. To this end, we  
11 simulate realistic multimodal time-series data for a prototypical remote worker,  
12 design a natural language prompt to elicit recommendations from an LLM, and  
13 employ an independent LLM to evaluate the generated outputs in terms of safety,  
14 relevance, and feasibility. The results suggest that LLMs are capable of detecting  
15 meaningful behavioral patterns and translating them into actionable guidance. This  
16 approach has the potential to support individuals in developing adaptive routines for  
17 health and productivity management. Future research should emphasize real-world  
18 validation, integration with digital health platforms, and the establishment of ethical  
19 safeguards.

20 **1 Motivation**

21 The organization of work has undergone profound transformations in recent decades. Increasingly  
22 flexible arrangements, hybrid work models, and the widespread adoption of remote working have  
23 blurred the boundaries between professional and private life. While these developments afford  
24 individuals greater autonomy, they simultaneously increase the demands on self-organization. Many  
25 workers report difficulties in maintaining structure, with consequences for sleep quality, physical  
26 activity, and mental health. Epidemiological data confirm substantial increases in the prevalence of  
27 burnout and depression over the past decade, with implications for individual well-being, workforce  
28 productivity, and public health systems (World Health Organization, 2022).

29 In conventional office settings, organizational routines such as fixed schedules, communal breaks, and  
30 social support structures provided external scaffolding for daily rhythms. The dissolution of such  
31 anchors in remote and hybrid contexts shifts responsibility for health and productivity management  
32 to individuals. However, effective self-management requires balancing multiple interacting factors,  
33 including sleep, stress, nutrition, physical activity, and social connectedness. Human introspection is  
34 often insufficient to capture these dynamic interrelations, leaving individuals unable to identify subtle  
35 but consequential patterns.

36 Advances in digital health technologies offer promising avenues for addressing this challenge.  
37 Wearables and smartphones now enable continuous monitoring of physiology and behavior, while

38 digital platforms facilitate real-time self-reporting. Nevertheless, a persistent gap remains between  
39 the availability of heterogeneous, temporally dense data streams and their translation into actionable  
40 insights. LLMs, with their capacity to process natural language and structured data, represent a potential  
41 solution. By integrating multimodal data, LLMs could generate personalized recommendations  
42 expressed in accessible, everyday language. This paper explores this proposition through a simulated  
43 case study.

## 44 **2 Related Work: Diary Studies and Personal Informatics**

45 Psychological and health research has a long tradition of employing diary methods to capture  
46 experiences in daily life. Approaches such as the Experience Sampling Method (ESM) and Ecological  
47 Momentary Assessment (EMA) allow participants to report moods, behaviors, and contextual factors  
48 in real time. These methodologies reduce recall bias and provide ecologically valid insights into  
49 fluctuating states [1]. Applications span clinical psychology, organizational behavior, and health  
50 research. For instance, Verhagen et al. [2] demonstrated the utility of EMA in psychiatric populations,  
51 while Fritz et al. [3] synthesized methodological best practices for intensive longitudinal designs.

52 Complementing subjective measures, digital phenotyping and personal sensing approaches leverage  
53 sensors embedded in smartphones and wearables to capture behavioral and physiological data  
54 unobtrusively. Mohr et al. [4] described personal sensing as a paradigm shift in mental health,  
55 enabling objective assessment of activity, sleep, and social interactions. Pizzoli et al. [5] reviewed  
56 advances in digital phenotyping, highlighting the utility of accelerometry, GPS, and app usage data as  
57 proxies for lifestyle behaviors. Torrado et al. [6], in the ActiveAgeing study, demonstrated the value  
58 of combining wearable sensors with qualitative methods to monitor older adults.

59 Although diary methods and personal informatics yield rich insights into daily life, interpretation  
60 typically remains the responsibility of researchers or individuals themselves. Ecological momentary  
61 interventions [7] have sought to deliver just-in-time feedback, yet the integration of multimodal data  
62 into automated, personalized guidance remains limited. This gap underscores the opportunity for  
63 LLMs to serve as intermediaries between raw data and actionable lifestyle advice.

## 64 **3 Analysis of the Potential of LLM-Based Lifestyle Adjustment 65 Recommendations**

### 66 **3.1 Multimodal Dataset**

67 To examine feasibility, we simulated a seven-day multimodal dataset for a fictional persona, Alice, a  
68 35-year-old remote software developer. Each day includes self-reported mood, perceived stress, and  
69 productivity (scales 1–10), alongside wearable-derived measures: daily step count, mean heart rate  
70 variability (HRV), sleep duration, sedentary hours, and self-logged social contacts (interactions >10  
71 minutes with colleagues, friends, or family).

72 The data indicate clear associations (see Tables 1 and 2): days with insufficient sleep (<6.5 h)  
73 correspond to lower mood, higher stress, reduced productivity, diminished HRV, fewer steps, and  
74 limited social interaction (Days 2, 3, and 6). Conversely, days with sufficient sleep, higher physical  
75 activity, and at least three social contacts (Days 1, 4, and 7) align with improved mood and productivity.  
76 These observations highlight the interdependence of sleep, activity, social engagement, and well-being.

### 77 **3.2 Prompt Design**

78 We developed a structured natural language prompt instructing an LLM to analyze Alice’s dataset  
79 and generate recommendations. The prompt emphasized safety, contextual relevance, and feasibility,  
80 requiring outputs to consist of complete sentences with concise justifications.

81 The following prompt (template) was used:

82 You are given a seven-day dataset from a 35-year-old remote software  
83 developer named Alice. The dataset includes daily values for mood,  
84 stress, productivity, steps, heart rate variability (HRV), sleep  
85 duration, sedentary hours, and social contacts.

Table 1: Summary of subjective measures averaged across stress-level groups

Group	Mood	Stress	Productivity
High stress (Days 2, 3, 6)	4.7	7.3	4.7
Moderate stress (Day 5)	6.0	6.0	6.0
Low stress (Days 1, 4, 7)	6.7	4.3	7.7

Table 2: Summary of objective measures averaged across stress-level groups

Group	Steps	HRV (ms)	Sleep (h)	Sed. (h)	Contacts
High stress (Days 2, 3, 6)	4,167	35.0	6.0	10.3	0.7
Moderate stress (Day 5)	6,100	39.0	6.9	9.0	2.0
Low stress (Days 1, 4, 7)	8,400	44.3	7.7	7.3	3.3

86 Here is the data: {DATA}  
 87 Your task is to:  
 88 1. Analyze the dataset as a whole, looking for patterns across the  
 89 seven days.  
 90 2. Generate lifestyle adjustment recommendations that are safe,  
 91 contextually relevant, and feasible for someone in Alice's  
 92 situation.  
 93 3. Present your recommendations in complete sentences.  
 94 4. Provide a brief justification for each recommendation,  
 95 explaining which data patterns support it.  
 96 5. Focus only on constructive, non-clinical advice (e.g., daily  
 97 routines, work breaks, activity balance, social habits, sleep  
 98 hygiene).  
 99 Avoid medical diagnoses or treatments. Write clearly and concisely.

### 100 3.3 Example Recommendations

101 The LLM examined the following data patterns:

- 102 **Mood, Stress, Productivity:** Mood and productivity are higher, and stress lower, on days  
 103 with more sleep, higher steps, and more social contacts (Days 4 and 7). Conversely, mood  
 104 and productivity drop when sleep is short, steps are low, and stress rises (Days 2–3, 6).
- 105 **Activity & Sedentary Time:** Days with fewer than 5,000 steps and long sedentary hours  
 106 (10–11 h) correspond to lower mood and productivity.
- 107 **Sleep:** Better days (4 and 7) align with 7.8–7.9 hours of sleep. Short sleep (5.9–6.2 h) is  
 108 linked with worse stress and productivity.
- 109 **HRV:** HRV is higher (45–46 ms) on good days, and lower (34–36 ms) on stressful, less  
 110 active days.
- 111 **Social Contacts:** More interactions (3–4 contacts) coincide with better mood and productiv-  
 112 ity.

113 It produced the following recommendations:

- 114 **Prioritize consistent sleep of 7–8 hours per night.** Justification: On days with longer sleep  
 115 (Days 4 and 7), Alice's mood, stress, and productivity improved. Maintaining this range  
 116 could help stabilize energy and performance.
- 117 **Incorporate short movement breaks during work hours to reduce sedentary time.**  
 118 Justification: Higher sedentary hours (10–11 h) were linked with lower steps, mood, and  
 119 productivity (Days 2–3, 6). Standing or stretching breaks every hour may help counteract  
 120 this.

Table 3: Evaluation of lifestyle recommendations on safety, relevance, and feasibility

Recommendation	Safety	Relevance	Feasibility
Prioritize consistent sleep of 7–8 hours per night	Safe; aligns with sleep guidelines	Data: longer sleep leads to better mood, stress, productivity	Realistic with good sleep hygiene
Incorporate short movement breaks during work hours	Safe; light activity, minimal risk	Data: high sedentary leads to lower mood/productivity	Highly feasible (standing, stretching, walking)
Aim for at least 7,000–8,000 steps most days	Safe; aligns with activity guidelines	Data: higher steps lead to better mood, HRV, productivity	Achievable, requires planning on busy days
Schedule intentional social interactions	Safe; psychologically beneficial	Data: more contacts lead to higher mood, productivity	Feasible, depends on schedule/social network
Plan focused work blocks with breaks	Safe; aligns with healthy work practices	Data: lower stress with higher productivity	Feasible with structured planning

- 121           3. **Aim for at least 7,000–8,000 steps most days.** Justification: Days with more steps (Day 4:  
 122           8,600; Day 7: 9,100) coincided with higher mood, productivity, and HRV. Increasing daily  
 123           activity appears beneficial.  
 124           4. **Schedule intentional social interactions throughout the week.** Justification: Social contact  
 125           was highest (3–4) on Alice’s best days for mood and productivity. Regular connections may  
 126           provide emotional balance and reduce stress.  
 127           5. **Plan focused work blocks with breaks to manage stress.** Justification: Stress was lowest  
 128           when productivity was high (Days 4 and 7). Structured routines that mix deep work and  
 129           recovery breaks could replicate these conditions more often.

### 130           3.4 Evaluation by a Secondary LLM

131           A secondary LLM evaluated the recommendations for safety, relevance, and feasibility (Table 3).  
 132           All were rated safe for a healthy adult, directly grounded in the observed patterns, and realistic  
 133           to implement. The evaluator further noted consistency with established psychological evidence,  
 134           particularly regarding the benefits of sleep, activity, and social support.  
 135           For evaluation, following prompt template was used:

136           You are given a set of lifestyle recommendations generated for a  
 137           healthy adult based on a seven-day dataset.  
 138           Your task is to evaluate each recommendation on three dimensions:  
 139           1. Safety - Is the recommendation safe for a healthy adult?  
 140           2. Relevance - Is it directly grounded in the observed data  
 141           patterns?  
 142           3. Feasibility - Is it realistic to implement in daily life?  
 143           Provide a brief judgment for each dimension, using complete  
 144           sentences. If applicable, note whether the recommendation is  
 145           consistent with established psychological or behavioral evidence  
 146           (e.g., benefits of sleep, activity, or social support).  
 147           Here are the recommendations: {RECOMMENDATIONS}

### 148           3.5 Discussion

149           We illustrated that LLMs can interpret multimodal datasets encompassing behavioral, physiological,  
 150           and social variables to generate contextually appropriate recommendations. Importantly, the design of  
 151           prompts proved critical in constraining outputs toward safe, specific, and actionable suggestions. The  
 152           inclusion of social interaction metrics expanded interpretive scope beyond physical and physiological

153 dimensions, underscoring the multidimensional nature of well-being. To extend beyond the proof-  
154 of-concept stage, future research should deploy the proposed workflow on real-world multimodal  
155 datasets collected through combinations of wearables and ecological momentary assessment (EMA)  
156 applications. Even small-scale pilot studies would provide valuable ecological validity, enabling  
157 comparison between LLM-derived recommendations and patterns observed in naturalistic settings.  
158 Such validation would directly address the limitations of simulation and strengthen the robustness of  
159 the approach.

## 160 **4 Limitations**

161 Several limitations of this exploratory work should be acknowledged. The analysis relied on a  
162 simulated dataset rather than empirical observations, which allowed for systematic illustration of the  
163 proposed workflow but restricted the ecological validity of the findings. Application to authentic  
164 multimodal data streams, such as those combining wearable sensing with ecological momentary  
165 assessment, will be necessary to establish robustness under naturalistic conditions. The evaluation of  
166 recommendations was conducted by a secondary LLM from the same model family as the generator.  
167 While this setup allowed scalable and efficient assessment, it risks bias through shared training data  
168 and architectural characteristics, meaning that the evaluator may replicate or reinforce the generator's  
169 assumptions rather than provide an independent appraisal. While this approach enabled rapid and  
170 scalable appraisal of safety, feasibility, and relevance, it cannot substitute for expert or clinical  
171 judgment, and recommendations deemed plausible by an LLM may nevertheless prove infeasible or  
172 ethically problematic in practice. The scope of the dataset was deliberately narrow, encompassing  
173 only a small set of daily variables over a seven-day period, whereas human health and behavior are  
174 shaped by a far broader constellation of contextual and longitudinal influences including nutrition,  
175 workload, socioeconomic status, and environment. Finally, the present design does not address  
176 questions of privacy, personalization boundaries, or user agency. The integration of LLM-generated  
177 recommendations into everyday life raises substantive ethical and governance challenges related to  
178 accountability, transparency, and informed consent, all of which must be addressed before real-world  
179 deployment can be responsibly pursued. Taken together, these constraints delineate clear priorities for  
180 future work, including empirical validation, expansion of multimodal scope, and the development  
181 of safeguards that enable the responsible integration of LLM-based recommendation systems into  
182 digital health practice.

## 183 **5 Conclusion and Outlook**

184 This study explored the potential of LLMs to generate lifestyle adjustment recommendations based  
185 on multimodal data integrating subjective and sensor-derived measures. Using a simulated dataset,  
186 we demonstrated that LLMs can detect meaningful behavioral patterns and transform them into  
187 actionable, safe, and feasible advice. The use of a secondary LLM as an evaluator further suggests  
188 pathways for embedding quality control mechanisms into recommendation pipelines.

189 More broadly, the dual-LLM pipeline—consisting of a generator model that produces candidate  
190 recommendations and an evaluator model that rates their safety, feasibility, and contextual fit—could  
191 be generalized as a reusable research workflow across domains. Beyond digital health, analogous  
192 “generate-then-evaluate” pipelines may support tasks in education, computational creativity, or  
193 human–computer interaction research, where iterative AI-based self-checking reduces the risk of  
194 unsafe or irrelevant outputs. Adapting the workflow to these domains would require tailoring  
195 evaluation criteria to context: for instance, pedagogical appropriateness and inclusivity in education,  
196 or novelty and coherence in creativity research. Safeguards would need to be domain-specific  
197 but consistently oriented toward ensuring that outputs remain relevant, non-harmful, and ethically  
198 appropriate.

199 Future research should validate these findings in real-world contexts, with naturalistic multimodal  
200 datasets and expert benchmarks. Even small-scale pilot studies—for example, involving 15–20  
201 participants over four to six weeks with combined wearable sensing and ecological momentary  
202 assessment (EMA) self-reports—could provide valuable ecological validity and enable systematic  
203 comparison between LLM-derived recommendations and observed behavioral patterns. Expanding  
204 the range of modalities beyond the current scope is also crucial. Nutrition, work environment factors  
205 (including ergonomics and screen time), socioeconomic context, and environmental exposures such

206 as light and noise represent particularly salient extensions, each with documented relevance for stress  
207 regulation, sleep, and well-being. Their inclusion would enhance personalization and ecological  
208 validity of recommendations.

209 At the same time, ethical considerations such as transparency, privacy, and the mitigation of over-  
210 reliance must be systematically addressed. Human-in-the-loop oversight can be structured so that  
211 experts shape the evaluation rubrics while end-users retain agency to accept, reject, or adapt AI  
212 suggestions. This layered approach maintains AI as the lead analytic agent but embeds essential  
213 checkpoints of human judgment. Procedurally, safeguards should include clear disclaimers that  
214 outputs are not medical advice, transparent documentation of data provenance, opt-in consent for data  
215 use, and privacy-preserving storage. Technically, accountability can be supported through audit trails  
216 and logging of evaluator decisions, allowing independent review of system behavior.

217 With these measures in place, embedding LLM-based recommendation systems into digital health  
218 and personal informatics platforms could provide substantial benefits in helping individuals navigate  
219 the demands of increasingly flexible and complex working environments.

220 **References**

- 221 [1] T. J. Trull and U. W. Ebner-Priemer. Using experience sampling methods/ecological momentary  
222 assessment (ESM/EMA) in clinical assessment and clinical research: introduction to the special  
223 section. *Psychological Assessment*, 21:457–462, 12 2009. doi: 10.1037/a0017653.
- 224 [2] S. J. Verhagen, L. Hasmi, M. Drukker, J. van Os, and P. A. Delespaul. Use of the experience  
225 sampling method in the context of clinical trials. *Evidence-Based Mental Health*, 19:86–89, 08  
226 2016. doi: 10.1136/ebmental-2016-102418.
- 227 [3] Jessica Fritz, Marilyn Piccirillo, Zachary Cohen, Madelyn Frumkin, Olivia Kirtley, Julia  
228 Moeller, Andreas Neubauer, Lesley Norris, Noémi Schuurman, Evelien Snippe, and Laura  
229 Bringmann. So you want to do esm? 10 essential topics for implementing the experience-  
230 sampling method. *Advances in Methods and Practices in Psychological Science*, 7, 09 2024. doi:  
231 10.1177/25152459241267912.
- 232 [4] David C. Mohr, Mi Zhang, and Stephen M. Schueller. Personal sensing: understanding mental  
233 health using ubiquitous sensors and machine learning. *Annual Review of Clinical Psychology*, 13:  
234 23–47, 05 2017. doi: 10.1146/annurev-clinpsy-032816-044949.
- 235 [5] S. F. M. Pizzoli, D. Monzani, L. Conti, G. Ferraris, R. Grasso, and G. Pravettoni. Issues  
236 and opportunities of digital phenotyping: ecological momentary assessment and behavioral  
237 sensing in protecting the young from suicide. *Frontiers in Psychology*, 14:1103703, 2023. doi:  
238 10.3389/fpsyg.2023.1103703.
- 239 [6] J. C. Torrado, B. S. Husebo, H. G. Allore, A. Erdal, S. E. Fæø, H. Reithe, et al. Digital phenotyping  
240 by wearable-driven artificial intelligence in older adults and people with parkinson's disease:  
241 Protocol of the mixed method, cyclic activeageing study. *PLOS ONE*, 17:e0275747, 10 2022.  
242 doi: 10.1371/journal.pone.0275747.
- 243 [7] Kristin E. Heron and Joshua M. Smyth. Ecological momentary interventions: incorporating  
244 mobile technology into psychosocial and health behaviour treatments. *British Journal of Health  
245 Psychology*, 15:1–39, 02 2010. doi: 10.1348/135910709X466063.

246 **Reproducibility Statement**

247 We have taken several steps to support reproducibility of this work. The study relies on a simulated  
248 seven-day multimodal dataset, which is illustrated in section 3 ("Analysis of the Potential of LLM-  
249 Based Lifestyle Adjustment Recommendations") and can be readily reconstructed by other researchers.  
250 All prompt designs and evaluation procedures are explicitly reported to enable replication. The large  
251 language model used was GPT-5, accessed through the standard web interface available to all users at  
252 the time of writing. No proprietary fine-tuning or hidden system settings were applied. While the  
253 exact outputs of generative models may vary slightly across runs, the reproducibility of the analysis  
254 lies in the transparent description of input data, prompts, and evaluation criteria.

255 The prompt used for generating the initial draft was as follows:

256 You are a researcher. Write a scientific paper in English language.  
257 Only use existing peer-reviewed scientific literature to write the  
258 paper. The length should be approximately 24,000 characters  
259 including spaces. Topic: Potential of LLM-generated lifestyle  
260 adaptation recommendations based on multimodal data

261 **0. Abstract and author information**

262 **1. Motivation (1 page)**

263 Content: Brief justification of why the topic is important,  
264 e.g., due to the increasing prevalence of burnouts and  
265 depression as well as flexibilization of work, more work from  
266 home (WFH), etc. All of this leads to the fact that people need  
267 to organize themselves better and actively manage their health,  
268 well-being, and productivity. Employers can no longer provide  
269 this, especially against the backdrop of working from home.

270 Overall message and conclusion: Self-management is becoming  
271 increasingly important, but it is difficult due to the many  
272 influencing variables. Therefore, IT support is important.  
273 **2. Related Work: Diary Studies and Personal Informatics (1 page)**  
274 Content: Personal Informatics and other similar forms such as  
275 diary studies have become established and are also researched  
276 under the title "Diary Study" or "Ecological Momentary  
277 Assessment" and "Ecological Momentary Intervention" in research,  
278 especially in organizational and health psychology. So far,  
279 however, everything is based on self-assessments, "objective"  
280 sensor data are used too little.  
281 Overall message and conclusion: Personal Informatics and diary  
282 studies, especially those with supplementary sensor data, have  
283 the potential to improve self-management by revealing  
284 favorable/unfavorable patterns that could not be identified  
285 through introspection and "a bit of self-reflection" alone.  
286 **3. Analysis of the potential of LLM-based lifestyle adaptation**  
287 **recommendations (4 pages)**  
288 Content:  
289     3.1 Dataset: Generation of exemplary multimodal time series data  
290         (self-assessment and sensor data)  
291     3.2 Generation of recommendations including mention of the LLM  
292         prompt  
293     3.3 Evaluation: Another LLM, which did not generate the  
294         recommendations, is instructed to evaluate the quality of the  
295         recommendations given (i) the information about a persona,  
296         (ii) their data, and (iii) the derived recommendations.  
297         Criteria could be safety, relevance, and ease of  
298         implementation of the recommendations.  
299 Overall message and conclusion: It is possible to derive  
300 lifestyle adaptation recommendations from an LLM based on time  
301 series data.  
302 **4. Conclusion and Outlook (0.5 page)**  
303 **5. References (1 page)**

304 **Agents4Science AI Involvement Checklist**

- 305 1. **Hypothesis development:** Hypothesis development includes the process by which you  
306 came to explore this research topic and research question. This can involve the background  
307 research performed by either researchers or by AI. This can also involve whether the idea  
308 was proposed by researchers or by AI.

309 Answer: **[A]**

310 Explanation: The research question and motivation were generated by the human authors,  
311 based on prior literature in digital health and LLM applications. AI was not involved in  
312 selecting the research problem or formulating the main hypotheses.

- 313 2. **Experimental design and implementation:** This category includes design of experiments  
314 that are used to test the hypotheses, coding and implementation of computational methods,  
315 and the execution of these experiments.

316 Answer: **[B]**

317 Explanation: Humans designed the simulated dataset structure and defined the workflow  
318 (prompting, recommendation generation, and secondary evaluation). AI tools were used  
319 to implement portions of the design, such as generating the synthetic data narrative and  
320 running the evaluation prompts.

- 321 3. **Analysis of data and interpretation of results:** This category encompasses any process to  
322 organize and process data for the experiments in the paper. It also includes interpretations of  
323 the results of the study.

324 Answer: **[D]**

325 Explanation: AI models produced the primary lifestyle recommendations from the simulated  
326 multimodal data and provided feasibility/safety assessments. Human authors reviewed these  
327 outputs but contributed minimal original analysis beyond framing the interpretation.

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

331 Answer: **[C]**

332 Explanation: The manuscript text was drafted in collaboration with AI. The AI provided  
333 structured sections, academic phrasing, and polished formatting. Human authors curated  
334 content, reorganized material, and ensured alignment with conference style and scientific  
335 standards.

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

338 Description: AI-generated outputs were coherent but sometimes overly generic, lacking  
339 nuanced awareness of context and ethical considerations. Models cannot independently  
340 validate feasibility, and their assessments may miss subtle clinical or methodological issues.  
341 Careful human oversight was required to ensure accuracy, relevance, and responsible framing.

342 **Agents4Science Paper Checklist**

343 **1. Claims**

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

346 Answer: [Yes]

347 Justification: The abstract and introduction describe the use of simulated multimodal data  
348 and large language models to generate lifestyle recommendations. They accurately reflect  
349 the scope (proof-of-concept exploration) without overstating generalizability.

350 Guidelines:

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

360 **2. Limitations**

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

362 Answer: [Yes]

363 Justification: A dedicated Limitations section (Section 4) highlights the use of simulated data,  
364 reliance on AI (GPT-5) as proxy evaluator, restricted variables, and ethical considerations  
365 such as privacy and clinical oversight.

366 Guidelines:

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

389 **3. Theory assumptions and proofs**

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

392 Answer: [NA]

393 Justification: The paper does not include formal theoretical results or proofs, as the work is  
394 exploratory and methodological in nature.

395 Guidelines:

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

#### 403 **4. Experimental result reproducibility**

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

407 Answer: [Yes]

408 Justification: The simulated dataset, prompting procedure, and evaluation design are fully  
409 described and illustrated in section 3 ("Analysis of the Potential of LLM-Based Lifestyle  
410 Adjustment Recommendations"). Although code is not provided, sufficient detail is available  
411 to reproduce the study design.

412 Guidelines:

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

#### 423 **5. Open access to data and code**

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

426 Answer: [No]

427 Justification: The paper does not release code or data, but this is not central to the contribution.  
428 The dataset is synthetic and fully described, allowing others to replicate without direct file  
429 release.

430 Guidelines:

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

#### 442 **6. Experimental setting/details**

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

446 Answer: [NA]

447 Justification: The study does not involve model training or hyperparameter optimization. Ex-  
448 perimental details consist of prompt design and evaluation procedures, which are documented  
449 in the paper.

450 Guidelines:

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

## 456 7. Experiment statistical significance

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

459 Answer: [NA]

460 Justification: The work is conceptual and uses a single simulated dataset with qualitative AI  
461 outputs. No statistical tests or error bars were applicable.

462 Guidelines:

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

## 470 8. Experiments compute resources

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

474 Answer: [NA]

475 Justification: No compute-intensive experiments were conducted. Only standard queries to  
476 an LLM chat interface (GPT-5 Web UI) were performed, requiring minimal resources.

477 Guidelines:

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

## 483 9. Code of ethics

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

486 Answer: [Yes]

487 Justification: The study uses only synthetic data, involves no human participants, and  
488 explicitly discusses ethical considerations such as privacy, accountability, and potential  
489 misuse.

490 Guidelines:

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

## 495 10. Broader impacts

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

498 Answer: [Yes]

499 Justification: The paper reflects on potential benefits (e.g., scalable digital health support)  
500 and risks (e.g., over-reliance on AI, privacy concerns, ethical governance), including  
501 considerations for mitigation.

502 Guidelines:

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