



PERSONA $\ddot{\text{X}}$: MULTIMODAL DATASETS WITH LLM-INFERRRED BEHAVIOR TRAITS

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding human behavior traits is central to applications in human-computer interaction, computational social science, and personalized AI systems. Such understanding often requires integrating multiple modalities to capture nuanced patterns and relationships. However, existing resources rarely provide datasets that combine behavioral descriptors with complementary modalities such as facial attributes and biographical information. To address this gap, we present **PersonaX**, a curated collection of multimodal datasets designed to enable comprehensive analysis of public human traits across modalities. **PersonaX** consists of (1) **CelebPersona**, featuring 9444 public figures from diverse occupations, and (2) **AthlePersona**, covering 4181 professional athletes across 7 major sports leagues. Each dataset includes behavioral trait assessments inferred by three high-performing large language models, alongside facial imagery and structured biographical features.

We analyze PersonaX at two complementary levels. First, we abstract high-level trait scores from text descriptions and apply five statistical independence tests to examine their relationships with other modalities. Second, we introduce a novel causal representation learning (CRL) framework tailored to multimodal and multi-measurement data, providing theoretical identifiability guarantees. Experiments on both synthetic and real-world data demonstrate the effectiveness of our approach. By unifying structured and unstructured analysis, PersonaX establishes a foundation for studying LLM-inferred behavioral traits in conjunction with visual and biographical attributes, advancing multimodal trait analysis and causal reasoning.

 CelebPersona: huggingface.co/datasets/Persona-X/celebpersona

 **AthlePersona:** huggingface.co/datasets/Persona-X/athlepersona

1 INTRODUCTION

Human behavioral traits (or behavior summaries) refer to interpretable patterns of conduct inferred from publicly available information, e.g., spoken or written language, facial expressions, and biographical records, that reflect how individuals present themselves to the world [??]. These traits are often related to, but not identical with, psychological notions of personality, which typically require self-reports or expert evaluation [??]. Unlike clinical diagnoses, public behavioral traits can be inferred ethically and at scale from observable information, providing reproducible, population-wide insights that complement traditional personality research without medicalizing individuals. In recent years, the feasibility of extracting or summarizing such traits has been greatly enhanced by the rapid growth of large language models (LLMs) [??]. Indeed, several studies suggest that LLM-based assessments of behavioral traits aligned with the Big Five framework can be valid and reliable under specific prompting configurations [??]. These LLM-driven approaches not only scale efficiently but may also mitigate biases inherent in self-reports.

Related Work. Several multimodal datasets have been developed to study behavioral or personality-related constructs. Notable examples include YouTube-Vlogs [2], FI-V2 [2], MuPTA [2], MDPE [2], and Amigos [2]. These datasets typically combine modalities such as video, audio, or physiological signals, and are primarily designed for prediction tasks like personality recognition, first-impression analysis, or deception detection. While valuable, most lack textual trait descriptions and offer no framework for interpreting or causally analyzing relationships across modalities. Complementing these resources, a

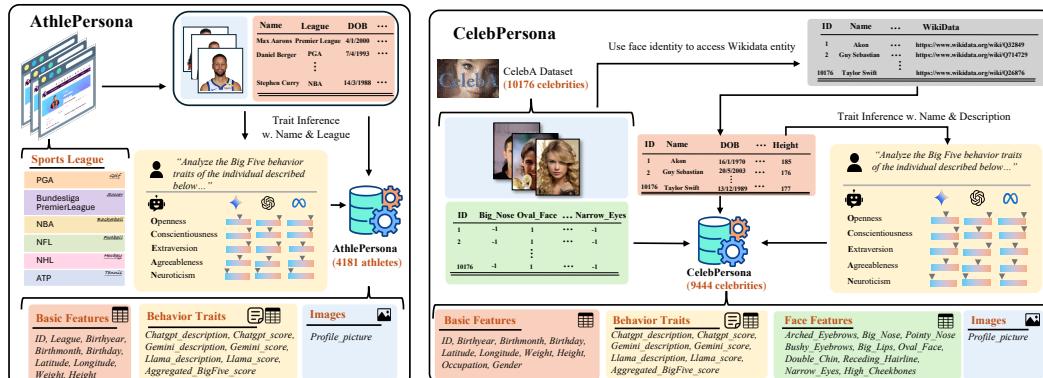


Figure 1: **Processing pipelines** to build AthlePersona (Left) and CelebPersona (Right) datasets. (1) AthlePersona was constructed by collecting player rosters and publicly available data (including facial images and basic features) from the official websites of prominent sports leagues. LLMs were then used for behavior traits inference. (2) CelebPersona was derived from CelebA. Face identities of celebrities are linked to their corresponding Wikidata entity, enabling the retrieval of additional biographical details and physical characteristics, and similarly processed through LLMs for behavior traits inference. All data modalities are finally combined together, followed by post-processing.

growing body of research has examined how observable features in one modality can signal traits in another. The “kernel of truth” hypothesis posits that physical characteristics may reflect underlying behavioral tendencies ?. Empirical studies support this view: facial features have been linked to health cues ?, aggression via facial width-to-height ratio ?, and trait judgments from body images ?. Advances in machine learning further show that Big Five traits can be predicted from static facial images with above-chance accuracy ??. Beyond vision, correlations have also been documented with activity levels ?, sensor data ?, and physiological measurements ?. Together, these findings underscore the importance of studying relationships across modalities rather than treating each in isolation. Yet existing datasets and methods remain limited for systematic cross-modal and causal analyses. More details on related work are in App. ??.

To mitigate these gaps, we introduce PersonaX, a curated collection of large-scale multimodal datasets that link LLM-inferred behavioral trait assessments with complementary modalities. These traits are systematically summarized from public information, including direct quotes from interviews, observed behaviors, career trajectories, and biographical details. For consistency and simplicity, we infer behavioral traits along the five dimensions of the Big Five framework ?, each with an associated score. Specifically, PersonaX comprises (i) CelebPersona, containing 9444 public figures from the CelebA dataset ?, and (ii) AthlePersona, with 4181 professional athletes across 7 major sports leagues. Each record integrates (1) textual trait descriptions and Big Five scores inferred by three high-performing LLMs, (2) facial images, and (3) structured biographical metadata. To safeguard privacy, we release only transformed embeddings rather than raw images or texts.

Our contributions are mainly twofold. (i) We release PersonaX, multimodal datasets that combine LLM-inferred behavioral traits, facial embeddings, and biographical metadata for large populations of public figures. (ii) We propose a two-level analysis framework: at the structured level, applying diverse statistical independence tests to uncover behavioral-trait dependencies; at the unstructured level, introducing a new causal representation learning approach with identifiability guarantees tailored for multimodal, multi-measurement settings. Experiments on synthetic and real-world data demonstrate the effectiveness of this framework. By bridging structured and unstructured analysis, PersonaX provides a resource for examining LLM-inferred behavioral traits alongside visual and biographical attributes, opening pathways for deeper multimodal trait interpretation and causal reasoning. Our long-term vision is to leverage such resources to uncover invariant causal patterns across populations, thereby fostering diversity, equality, and mutual respect for all human beings.

108
 109
 110
 111
 112
Table 1: Comparative evaluations on some of the state-of-the-art LLMs using CelebPersona and
 AthlePersona subsets. Metrics consist of generation time (GT), missing rate (MR), indecisive rate
 (IR), privacy preservation (PP), output formatting (OF), context consistency (CC), factual accuracy
 (FA), and an overall score (OS). Please refer to Sec. ?? and App. ?? for more details.

Model (LLMs)	CelebPersona								AthlePersona							
	GT↓	MR↓	IR↓	PP↑	OF↑	CC↑	FA↑	OS↑	GT↓	MR↓	IR↓	PP↑	OF↑	CC↑	FA↑	OS↑
ChatGPT-4o-Latest	4.19	0.03	0.17	0.99	1.00	1.00	1.00	0.96	3.92	0.27	0.17	1.00	1.00	0.99	1.00	0.93
Gemini-2.5-Pro	23.48	0.06	0.19	0.99	1.00	1.00	1.00	0.96	21.31	0.29	0.22	0.99	1.00	1.00	1.00	0.91
Qwen2.5-Max	9.10	0.24	0.29	1.00	0.99	0.99	1.00	0.91	8.93	0.32	0.36	1.00	0.99	0.99	1.00	0.88
Grok-3-Beta	5.92	0.34	0.17	1.00	1.00	0.99	1.00	0.91	4.96	0.66	0.10	1.00	1.00	1.00	1.00	0.87
Llama-4-Maverick	3.73	0.25	0.29	1.00	1.00	0.97	1.00	0.90	3.99	0.30	0.43	1.00	1.00	0.95	1.00	0.87
Gemini-2.0-Flash-T.	8.83	0.28	0.38	0.99	1.00	1.00	1.00	0.89	8.26	0.48	0.27	0.97	1.00	0.99	1.00	0.87
DeepSeek-R1	39.13	0.40	0.11	0.98	0.90	1.00	1.00	0.89	26.61	0.64	0.10	1.00	0.81	1.00	0.98	0.84
Qwen-Plus	9.22	0.24	0.46	1.00	0.99	1.00	1.00	0.88	8.87	0.30	0.48	1.00	0.98	1.00	1.00	0.87
DeepSeek-V3-0324	14.18	0.47	0.24	0.99	1.00	1.00	1.00	0.88	7.75	0.64	0.20	1.00	1.00	1.00	1.00	0.86
Gemini-2.0-Flash	2.53	0.43	0.32	1.00	1.00	0.99	1.00	0.87	2.29	0.69	0.18	1.00	1.00	1.00	1.00	0.86

2 PERSONAX DATASET

123
 124
AthlePersona. This dataset was built from scratch, documenting 4181 male professional athletes
 125 across 7 major sport leagues worldwide, including the NBA, NFL, NHL, ATP, PGA, Premier League,
 126 and Bundesliga. From official league sources, we collected biographical information (e.g., name,
 127 birth date, nationality), physical attributes (e.g., height, weight), and facial images. Nationalities were
 128 geocoded into continuous spatial coordinates (latitude and longitude) to support geographic analysis.
 129

130
CelebPersona. This dataset was built upon the established CelebA dataset ? with rich facial
 131 attribute annotations. We linked each face identity to its corresponding WikiData entity, enabling
 132 retrieval of additional biographical details and physical characteristics. From the original 40 CelebA
 133 attributes, we manually retained 10 (e.g., *Big Nose*, *High Cheekbones*) that are more stable and likely
 134 to reflect inherent appearance properties, while discarding those subject to short-term variation (e.g.,
 135 *Heavy Makeup*). In total, our dataset contains 9444 public figures across diverse occupations.
 136

137 **Multimodality.** Each record integrates three components: (1) textual behavioral trait descriptions
 138 and Big Five scores inferred by LLMs, (2) facial images or embeddings with annotations, and (3)
 139 structured biographical metadata. Dataset statistics, distributions, and missingness are in App. ??.

2.1 LLM SELECTION AND PROMPT DESIGN FOR BEHAVIOR TRAIT INFERENCE

140 We systematically evaluated Top eight leading LLMs from the Arena leaderboard ? and two additional
 141 strong performers (Qwen2.5-Max, QwQ-32B). Models were assessed on eight criteria, including
 142 generation time, missing/indecisive rates, privacy preservation, factual accuracy, context consistency,
 143 and more (see App. ??). Prompts were carefully designed to balance interpretability and consistency.
 144 We experimented with numeric/textual outputs, 3-level/5-level scoring scales, and ordering directions,
 145 running controlled trials across models. Results showed that coarse numeric 3-point scales minimized
 146 variability, while more complex textual scales increased inconsistency (see App. ??). Based on these
 147 evaluations, we ultimately selected three consistently best-performing LLMs (ChatGPT-4o-Latest,
 148 Gemini-2.5-Pro, and Llama-4-Maverick) to generate trait descriptions and corresponding scores. See
 149 Prompt ?? for full prompt. For more experiments refer to App. ??.
 150

2.2 ETHICAL CONSIDERATIONS: CONSENT, PRIVACY, BIAS, AND USAGE

151 To address ethical and technical concerns, we emphasize four aspects: (i) *Consent and legality*:
 152 both datasets are derived entirely from legally accessible, consent-based resources (see Tab. ??),
 153 including the seven official sport league websites, CelebA (non-commercial academic use), and
 154 WikiData (free for unrestricted use). All data are collected strictly for non-commercial, academic
 155 research; (ii) *Privacy protection*: no raw images or textual descriptions are released. Each facial
 156 image is replaced with a 1024-dimensional embedding and each textual trait description with a 3584-
 157 dimensional embedding, further obfuscated through an additional invertible transformation, while
 158 categorical variables are converted into indices; (iii) *Bias and limitations*: AthlePersona includes
 159 only male athletes, while CelebPersona reflects wealthy, high-visibility individuals. Findings
 160

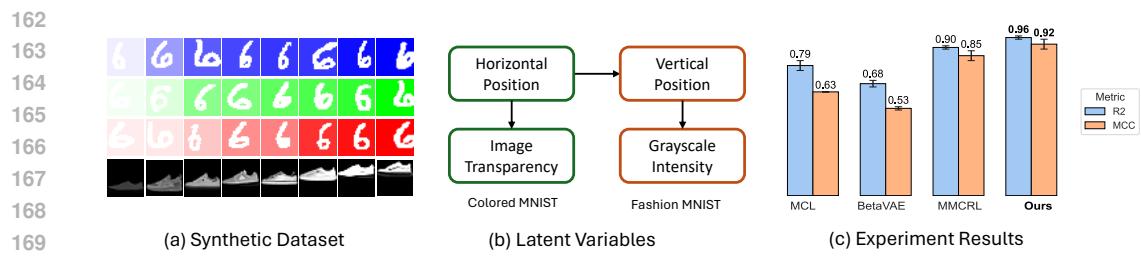


Figure 2: **Synthetic experiments.** (a) Synthetic Colored MNIST and fashion MNIST modalities. (b) Ground-truth causal graph. (c) Results: our method outperforms baselines in R^2 and MCC.

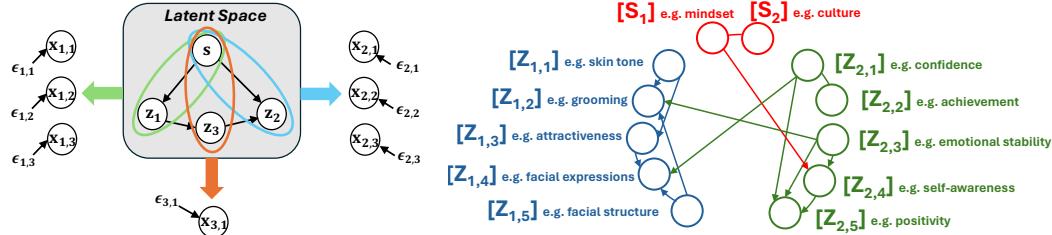


Figure 3: **Causal modeling.** Left: multimodal, multi-measurement ground-truth causal model. Right: estimated latent causal graph from AthlePersona, presenting interpretable causal relations.

should therefore be interpreted as population-specific rather than universal, though the diversity across domains offers opportunities to uncover invariant causal patterns; and (iv) *Usage restrictions*: a mandatory `USAGE_GUIDELINES.md` limits the dataset to non-commercial use and requires users to refrain from applying it (or any derived models) in high-stakes contexts (e.g., employment, insurance, or lending).

3 TWO-LAYER CAUSAL ANALYSIS FRAMEWORK

We analyze PersonaX at two complementary levels. **(i) Structured dependence discovery.** For each individual, three LLMs generate textual trait descriptions that are abstracted into Big Five scores and aggregated; in CelebPersona, multiple image attributes are also aggregated. We then apply five independence tests, KCI ?, RCIT ?, HSIC ?, Chi-square ?, and G-square ?, to measure dependencies between trait scores and tabular features ($p < 0.05$). Results in Tab.?? reveal distinct patterns: in CelebPersona, demographic and facial attributes exhibit broad dependence with traits, while in AthlePersona, league affiliation, birth year, and physical attributes dominate. Geographic factors show consistent, moderate effects across both datasets. These results highlight population-specific transfer mechanisms, with appearance cues stronger for public figures and organizational context stronger for athletes (see App.??). **(ii) Unstructured causal representation learning.** We further learn latent variables and their causal relations directly from text and images via a multimodal, multi-measurement CRL framework. Fig.??(Left) shows the causal modelling. Shared latents (s) capture cross-modal structure, while modality-specific latents (z) represent unique factors. Multiple measurements (x), e.g., multiple LLM assessments, provide variability essential for identifiability, which we establish theoretically (see App. ?? and ??). **Experiments** on synthetic MNIST variants validate our method: the performance comparisons in Fig. ?? show our approach outperforming BetaVAE ?, MCL ?, and MMCRL ? in both R^2 and MCC. On real data, the estimated causal graph from AthlePersona in Fig. ??(Right) aligns well with the assumed multimodal generative structure, revealing interpretable cross-modal pathways between text- and image-derived latents. Results for CelebPersona are provided in Fig.?. Please refer to App. ?? and ?? for more details.

4 CONCLUSION

We presented PersonaX, two multimodal datasets linking LLM-inferred behavioral traits with facial and biographical information. Our two-level analysis pipeline combines structured dependence tests with unstructured causal representation learning, addressing both theoretical and empirical

216 aspects: from the theoretical side, we propose the novel identifiability theory tailored for multimodal,
 217 multi-measurement CRL; empirically, we demonstrate population-specific patterns and interpretable
 218 latent structures. These resources provide a foundation for studying invariant causal mechanisms of
 219 human behavioral traits while promoting diversity, equality, and mutual respect for all human beings.
 220

221 REFERENCES

223 Atpt terms of use. <https://www.atptour.com/en/terms-and-conditions>. Accessed:
 224 Aug 2025.

225 Bundesliga terms of use. <https://www.bundesliga.com/en/bundesliga/info/terms-of-use-services>. Accessed: Aug 2025.

228 Laliga terms of use. <https://www.laliga.com/en-GB/legal/legal-web>. Accessed:
 229 Aug 2025.

230 Legaseriea terms of use. https://img.legaseriea.it/vimages/64ca8e48/INTERNATIONAL%20MEDIA%20RIGHTS_GENERAL%20TERMS%20AND%20CONDITIONS%20OF%20THE%20LICENSE%20AGREEMENT.pdf. Accessed: Aug 2025.

234 Ligue1 terms of use. <https://ligue1.com/en/legal/cgu>. Accessed: Aug 2025.

236 Mlb terms of use. <https://www.mlb.com/official-information/terms-of-use>. Accessed: Aug 2025.

238 Nba terms of use. <https://www.nba.com/termsofuse>. Accessed: Aug 2025.

240 Nfl terms of use. <https://www.nfl.com/legal/terms/>. Accessed: Aug 2025.

241 Nhl terms of use. <https://www.nhl.com/info/terms-of-service>. Accessed: Aug 2025.

244 Pga terms of use. <https://www.pgatour.com/company/terms-of-use>. Accessed: Aug 2025.

246 Premierleague terms of use. <https://www.premierleague.com/en/terms-and-conditions>. Accessed: Aug 2025.

249 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 250 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 251 *arXiv preprint arXiv:2303.08774*, 2023.

254 Kartik Ahuja, Amin Mansouri, and Yixin Wang. Multi-domain causal representation learning via
 253 weak distributional invariances, 2023. URL <https://arxiv.org/abs/2310.02854>.

255 Anthropic. Introducing claudie, 2023. URL <https://www.anthropic.com/index/introducing-claudie>.

257 Syeda Asra and D. Shubhangi. Personality trait identification using unconstrained cursive and mood
 258 invariant handwritten text. *International Journal of Education and Management Engineering*, 5:
 259 20–31, 10 2015. doi: 10.5815/ijeme.2015.05.03.

260 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,
 261 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

263 Joan-Isaac Biel and Daniel Gatica-Perez. The youtube lens: Crowdsourced personality impressions
 264 and audiovisual analysis of vlogs. *IEEE Transactions on Multimedia*, 15(1):41–55, 2012.

265 Stephen R Briggs and Jonathan M Cheek. The role of factor analysis in the development and
 266 evaluation of personality scales. *Journal of personality*, 54(1):106–148, 1986.

268 Cong Cai, Shan Liang, Xuefei Liu, Kang Zhu, Zhengqi Wen, Jianhua Tao, Heng Xie, Jizhou Cui,
 269 Yiming Ma, Zhenhua Cheng, et al. Mdpe: A multimodal deception dataset with personality and
 emotional characteristics. *arXiv preprint arXiv:2407.12274*, 2024.

- 270 Justin M Carré and Cheryl M McCormick. In your face: facial metrics predict aggressive behaviour
 271 in the laboratory and in varsity and professional hockey players. *Proceedings of the Royal Society*
 272 *B: Biological Sciences*, 275(1651):2651–2656, 2008.
- 273
- 274 Raymond B. Cattell, Herbert W. Eber, and Maurice M. Tatsuoka. *Personality and Mood by Question-*
 275 *nnaire*. Institute for Personality and Ability Testing, 1970.
- 276 Fabio Celli, Elia Bruni, and Bruno Lepri. Automatic personality and interaction style recognition
 277 from facebook profile pictures. In *Proceedings of the 22nd ACM International Conference on*
 278 *Multimedia*, MM ’14, pp. 1101–1104, New York, NY, USA, 2014. Association for Computing
 279 Machinery. ISBN 9781450330633. doi: 10.1145/2647868.2654977. URL <https://doi.org/10.1145/2647868.2654977>.
- 280
- 281 Mark Chen, Jerry Tworek, Heewoo Jun, et al. Evaluating large language models trained on code. In
 282 *arXiv preprint arXiv:2107.03374*, 2021.
- 283
- 284 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng
 285 Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open
 286 platform for evaluating llms by human preference. In *Forty-first International Conference on*
 287 *Machine Learning*, 2024.
- 288
- 289 David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine*
 learning research, 3(Nov):507–554, 2002.
- 290
- 291 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, et al. Palm: Scaling language modeling with
 292 pathways. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2022.
- 293
- 294 Deborah A Cobb-Clark and Stefanie Schurer. The stability of big-five personality traits. *Economics*
 Letters, 115(1):11–15, 2012.
- 295
- 296 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, et al. Training verifiers to solve math word
 297 problems. In *arXiv preprint arXiv:2110.14168*, 2021.
- 298
- 299 John B Conway. *A course in functional analysis*, volume 96. Springer Science & Business Media,
 1994.
- 300
- 301 James Cussens. Bayesian network learning with cutting planes. In *Proceedings of the Twenty-Seventh*
 302 *Conference on Uncertainty in Artificial Intelligence*, pp. 153–160, 2011.
- 303
- 304 Imant Daunhawer, Alice Bizeul, Emanuele Palumbo, Alexander Marx, and Julia E Vogt. Identifiability
 305 results for multimodal contrastive learning. In *The Eleventh International Conference on Learning*
 Representations, 2023.
- 306
- 307 Kristina M DeNeve and Harris Cooper. The happy personality: a meta-analysis of 137 personality
 308 traits and subjective well-being. *Psychological bulletin*, 124(2):197, 1998.
- 309
- 310 Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xiang Liu, and Tong Zhang. Active prompting
 311 with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*, 2023.
- 312
- 313 Dario Dotti, Mirela Popa, and Stylianos Asteriadis. Behavior and personality analysis in a nonso-
 314 cial context dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
 Recognition Workshops, pp. 2354–2362, 2018.
- 315
- Nelson Dunford and Jacob T. Schwartz. *Linear Operators*. John Wiley & Sons, New York, 1971.
- 316
- Hugo Jair Escalante, Heysem Kaya, Albert Ali Salah, Sergio Escalera, Yağmur Güçlütürk, Umut
 317 Güçlü, Xavier Baró, Isabelle Guyon, Julio CS Jacques Junior, Meysam Madadi, et al. Modeling,
 318 recognizing, and explaining apparent personality from videos. *IEEE Transactions on Affective*
 319 *Computing*, 13(2):894–911, 2020.
- 320
- Hans J. Eysenck and Sybil B. G. Eysenck. *Manual of the Eysenck Personality Questionnaire*. Hodder
 321 and Stoughton, 1975.
- 322
- Shunxing Fan, Mingming Gong, and Kun Zhang. On the recoverability of causal relations from
 323 temporally aggregated iid data. *arXiv preprint arXiv:2406.02191*, 2024.

- 324 Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds*
 325 and *Machines*, 30:681–694, 2020.
- 326
- 327 Nan Gao, Wei Shao, and Flora D Salim. Predicting personality traits from physical activity intensity.
 328 *Computer*, 52(7):47–56, 2019.
- 329
- 330 Alan S Gerber, Gregory A Huber, David Doherty, and Conor M Dowling. The big five personality
 331 traits in the political arena. *Annual Review of Political Science*, 14(1):265–287, 2011.
- 332 Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand
 333 Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the*
 334 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 15180–15190, 2023.
- 335
- 336 Lewis R. Goldberg. The structure of phenotypic personality traits. *American Psychologist*, 48(1):
 337 26–34, 1993.
- 338 Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical
 339 dependence with hilbert-schmidt norms. In *International conference on algorithmic learning*
 340 *theory*, pp. 63–77. Springer, 2005.
- 341
- 342 Yağmur Güçlütürk, Umut Güçlü, Xavier Baró, Hugo Jair Escalante, Isabelle Guyon, Sergio Escalera,
 343 Marcel A.J. van Gerven, and Rob van Lier. Multimodal first impression analysis with deep residual
 344 networks. *IEEE Transactions on Affective Computing*, 9(3):316–329, 2018. doi: 10.1109/TAFFC.
 345 2017.2751469.
- 346 Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botvinick,
 347 Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a
 348 constrained variational framework. *ICLR (Poster)*, 3, 2017.
- 349
- 350 Biwei Huang, Kun Zhang, Jiji Zhang, Joseph Ramsey, Ruben Sanchez-Romero, Clark Glymour,
 351 and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *Journal of*
 352 *Machine Learning Research*, 21(89):1–53, 2020.
- 353 Chin-Wei Huang, David Krueger, Alexandre Lacoste, and Aaron Courville. Neural autoregressive
 354 flows. In *International conference on machine learning*, pp. 2078–2087. PMLR, 2018.
- 355
- 356 Aapo Hyvärinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ica using auxiliary variables and
 357 generalized contrastive learning. In *The 22nd International Conference on Artificial Intelligence*
 358 *and Statistics*, pp. 859–868. PMLR, 2019.
- 359 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
 360 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
 361 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 362
- 363 Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. Personallm:
 364 Investigating the ability of large language models to express personality traits. *arXiv preprint*
 365 *arXiv:2305.02547*, 2023.
- 366 John A Johnson. Units of analysis for the description and explanation of personality. In *Handbook of*
 367 *personality psychology*, pp. 73–93. Elsevier, 1997.
- 368
- 369 Alexander Kachur, Evgeny Osin, Denis Davydov, Konstantin Shutilov, and Alexey Novokshonov.
 370 Assessing the big five personality traits using real-life static facial images. *Scientific Reports*, 10
 371 (1):8487, 2020.
- 372 T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial net-
 373 works. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 374 pp. 4401–4410, 2019.
- 375
- 376 Ilyes Khemakhem, Diederik Kingma, Ricardo Monti, and Aapo Hyvärinen. Variational autoencoders
 377 and nonlinear ica: A unifying framework. In *International conference on artificial intelligence*
 378 *and statistics*, pp. 2207–2217. PMLR, 2020.

- 378 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
 379 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural*
 380 *information processing systems*, 33:18661–18673, 2020.
- 381
- 382 Mikko Koivisto and Kismat Sood. Exact bayesian structure discovery in bayesian networks. *The*
 383 *Journal of Machine Learning Research*, 5:549–573, 2004.
- 384
- 385 Meera Komaraju, Steven J Karau, Ronald R Schmeck, and Alen Avdic. The big five personality
 386 traits, learning styles, and academic achievement. *Personality and individual differences*, 51(4):
 387 472–477, 2011.
- 388
- 389 M. Kosinski, D. Stillwell, and T. Graepel. Facebook as a research tool for the social sciences:
 390 Opportunities, challenges, ethical considerations, and practical guidelines. *American Psychologist*,
 391 70(6):543–556, 2015.
- 392
- 393 Robin SS Kramer and Robert Ward. Internal facial features are signals of personality and health.
Quarterly Journal of Experimental Psychology, 63(11):2273–2287, 2010.
- 394
- 395 K. Kärkkäinen and J. Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias
 396 measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications*
 397 *of Computer Vision*, pp. 1548–1558, 2021.
- 398
- 399 Loka Li, Zhenhao Chen, Guangyi Chen, Yixuan Zhang, Yusheng Su, Eric Xing, and Kun Zhang.
 400 Confidence matters: Revisiting intrinsic self-correction capabilities of large language models.
arXiv preprint arXiv:2402.12563, 2024a.
- 401
- 402 Loka Li, Haoyue Dai, Hanin Al Ghothani, Biwei Huang, Jiji Zhang, Shahar Harel, Isaac Bentwich,
 403 Guangyi Chen, and Kun Zhang. On causal discovery in the presence of deterministic relations. In
The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024b.
- 404
- 405 Loka Li, Ignavier Ng, Gongxu Luo, Biwei Huang, Guangyi Chen, Tongliang Liu, Bin Gu, and Kun
 406 Zhang. Federated causal discovery from heterogeneous data. *arXiv preprint arXiv:2402.13241*,
 407 2024c.
- 408
- 409 Longkang Li and Baoyuan Wu. Learning to accelerate approximate methods for solving integer
 410 programming via early fixing. *arXiv preprint arXiv:2207.02087*, 2022.
- 411
- 412 Longkang Li, Xiaojin Fu, Hui-Ling Zhen, Mingxuan Yuan, Jun Wang, Jiawen Lu, Xialiang Tong,
 413 Jia Zeng, and Dirk Schnieders. Bilevel learning for large-scale flexible flow shop scheduling.
Computers & Industrial Engineering, 168:108140, 2022.
- 414
- 415 Longkang Li, Siyuan Liang, Zihao Zhu, Chris Ding, Hongyuan Zha, and Baoyuan Wu. Learning to
 416 optimize permutation flow shop scheduling via graph-based imitation learning. In *Proceedings of*
the AAAI Conference on Artificial Intelligence, volume 38, pp. 20185–20193, 2024d.
- 417
- 418 Zijian Li, Shunxing Fan, Yujia Zheng, Ignavier Ng, Shaoan Xie, Guangyi Chen, Xinshuai Dong,
 419 Ruichu Cai, and Kun Zhang. Synergy between sufficient changes and sparse mixing procedure for
 420 disentangled representation learning. *arXiv preprint arXiv:2503.00639*, 2025.
- 421
- 422 Juan Lin. Factorizing multivariate function classes. *Advances in neural information processing*
 423 *systems*, 10, 1997.
- 424
- 425 Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard De Melo, Xiaogang Wang, Jifeng Dai,
 426 Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *European*
Conference on Computer Vision, pp. 388–404. Springer, 2022.
- 427
- 428 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
 429 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*
arXiv:2412.19437, 2024.
- 430
- 431 Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of*
the IEEE International Conference on Computer Vision, pp. 3730–3738, 2015.

- 432 Gongxu Luo, Haoyue Dai, Boyang Sun, Loka Li, Biwei Huang, Petar Stojanov, and Kun Zhang.
 433 Gene regulatory network inference in the presence of selection bias and latent confounders. *arXiv*
 434 *preprint arXiv:2501.10124*, 2025.
- 435 Muhammad Arslan Manzoor, Sarah Albarri, Ziting Xian, Zaiqiao Meng, Preslav Nakov, and Shang-
 436 song Liang. Multimodality representation learning: A survey on evolution, pretraining and its
 437 applications. *ACM Transactions on Multimedia Computing, Communications and Applications*, 20
 438 (3):1–34, 2023.
- 439 Haiyi Mao, Hongfu Liu, Jason Xiaotian Dou, and Panayiotis V Benos. Towards cross-modal causal
 440 structure and representation learning. In *Machine Learning for Health*, pp. 120–140. PMLR, 2022.
- 441 Juan Abdon Miranda-Correa, Mojtaba Khomami Abadi, Nicu Sebe, and Ioannis Patras. Amigos: A
 442 dataset for affect, personality and mood research on individuals and groups. *IEEE transactions on*
 443 *affective computing*, 12(2):479–493, 2018.
- 444 Gelareh Mohammadi, Alessandro Vinciarelli, and Marcello Mortillaro. The voice of personality:
 445 Mapping nonverbal vocal behavior into trait attributions. *Proceedings of the 2nd international*
 446 *workshop on Social signal processing. ACM*, 10 2010. doi: 10.1145/1878116.1878123.
- 447 Isabel Briggs Myers, Mary H. McCaulley, Naomi L. Quenk, and Allen L. Hammer. *MBTI Manual: A*
 448 *guide to the development and use of the Myers-Briggs Type Indicator*. Consulting Psychologists
 449 Press, Palo Alto, CA, 3rd edition, 1998.
- 450 Laura P Naumann, Simine Vazire, Peter J Rentfrow, and Samuel D Gosling. Personality judgments
 451 based on physical appearance. *Personality and social psychology bulletin*, 35(12):1661–1671,
 452 2009.
- 453 Ignavier Ng, AmirEmad Ghassami, and Kun Zhang. On the role of sparsity and dag constraints for
 454 learning linear dags. *Advances in Neural Information Processing Systems*, 33:17943–17954, 2020.
- 455 Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive
 456 coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 457 OpenAI. Gpt-4 technical report, 2023. URL <https://openai.com/research/gpt-4>.
- 458 Atsushi Oshio, Kanako Taku, Mari Hirano, and Gul Saeed. Resilience and big five personality traits:
 459 A meta-analysis. *Personality and individual differences*, 127:54–60, 2018.
- 460 Ashwin Paranjape et al. Art: Automatic reasoning and tool-use for large language models. In
 461 *International Conference on Learning Representations (ICLR)*, 2023.
- 462 G. Park, H.A. Schwartz, J.C. Eichstaedt, M.L. Kern, M. Kosinski, D.J. Stillwell, L.H. Ungar, and
 463 M.E. Seligman. Automatic personality assessment through social media language. *Journal of*
 464 *Personality and Social Psychology*, 108(6):934–952, 2015.
- 465 Xiaokang Peng, Yake Wei, Andong Deng, Dong Wang, and Di Hu. Balanced multimodal learning via
 466 on-the-fly gradient modulation. In *Proceedings of the IEEE/CVF conference on computer vision*
 467 *and pattern recognition*, pp. 8238–8247, 2022.
- 468 James Pennebaker, Martha Francis, and Roger Booth. Linguistic inquiry and word count (liwc). 01
 469 1999.
- 470 Heinrich Peters, Moran Cerf, and Sandra C. Matz. Large language models can infer personality from
 471 free-form user interactions. *arXiv preprint arXiv:2405.13052*, 2024.
- 472 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 473 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 474 models from natural language supervision. In *International conference on machine learning*, pp.
 475 8748–8763. PMLR, 2021.
- 476 Haocong Rao, Cyril Leung, and Chunyan Miao. Can chatgpt assess human personalities? a general
 477 evaluation framework, 2023. URL <https://arxiv.org/abs/2303.01248>.

- 486 Eric Rawls, Erich Kummerfeld, and Anna Zilverstand. An integrated multimodal model of alcohol
 487 use disorder generated by data-driven causal discovery analysis. *Communications biology*, 4(1):
 488 435, 2021.
- 489
- 490 Sonia Roccas, Lilach Sagiv, Shalom H Schwartz, and Ariel Knafo. The big five personality factors
 491 and personal values. *Personality and social psychology bulletin*, 28(6):789–801, 2002.
- 492 J Peter Rothe. *The scientific analysis of personality*. Routledge, 2017.
- 493
- 494 Elena Ryumina, Dmitry Ryumin, Maxim Markitantov, Heysem Kaya, Alexey Karpov, et al. Multi-
 495 modal personality traits assessment (mupta) corpus: The impact of spontaneous and read speech.
 496 In *Proceedings of ISCA International Conference INTERSPEECH*, pp. 4049–4053, 2023.
- 497
- 498 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,
 499 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the
 IEEE*, 109(5):612–634, 2021.
- 500
- 501 Gregory Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Marwa Abdul-
 502 hai, Aleksandra Faust, and Maja Matarić. Personality traits in large language models. 2023.
- 503
- 504 Peter Spirtes and Clark Glymour. An algorithm for fast recovery of sparse causal graphs. *Social
 science computer review*, 9(1):62–72, 1991.
- 505
- 506 Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, prediction, and search*. MIT press,
 507 2001.
- 508
- 509 Peter L Spirtes, Christopher Meek, and Thomas S Richardson. Causal inference in the presence of
 510 latent variables and selection bias. *Conference on Uncertainty in Artificial Intelligence*, 1995.
- 511
- 512 Aarohi Srivastava et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of
 513 language models. In *arXiv preprint arXiv:2206.04615*, 2022.
- 514
- 515 Eric V Strobl, Kun Zhang, and Shyam Visweswaran. Approximate kernel-based conditional indepen-
 516 dence tests for fast non-parametric causal discovery. *Journal of Causal Inference*, 7(1):20180017,
 517 2019.
- 518
- 519 Nils Sturma, Chandler Squires, Mathias Drton, and Caroline Uhler. Unpaired multi-domain causal
 520 representation learning. *Advances in Neural Information Processing Systems*, 36, 2023.
- 521
- 522 Yuewen Sun, Lingjing Kong, Guangyi Chen, Loka Li, Gongxu Luo, Zijian Li, Yixuan Zhang, Yujia
 523 Zheng, Mengyue Yang, Petar Stojanov, et al. Causal representation learning from multimodal
 524 biomedical observations. In *The Thirteenth International Conference on Learning Representations*,
 525 2025.
- 526
- 527 Zeyi Sun, Ye Fang, Tong Wu, Pan Zhang, Yuhang Zang, Shu Kong, Yuanjun Xiong, Dahua Lin, and
 528 Jiaqi Wang. Alpha-clip: A clip model focusing on wherever you want. In *Proceedings of the
 IEEE/CVF conference on computer vision and pattern recognition*, pp. 13019–13029, 2024.
- 529
- 530 Ronald J Tallarida, Rodney B Murray, Ronald J Tallarida, and Rodney B Murray. Chi-square test.
 531 *Manual of pharmacologic calculations: with computer programs*, pp. 140–142, 1987.
- 532
- 533 Zeyu Tang, Zhenhao Chen, Loka Li, Xiangchen Song, Yunlong Deng, Yifan Shen, Guangyi Chen,
 534 Peter Spirtes, and Kun Zhang. Reflection-window decoding: Text generation with selective
 535 refinement. *arXiv preprint arXiv:2502.03678*, 2025.
- 536
- 537 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
 538 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
 539 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 540
- 541 Marina Tiuleneva, Vadim A. Porvatov, and Carlo Strapparava. Big-five backstage: A dramatic
 542 dataset for characters personality traits & gender analysis. In Michael Zock, Emmanuele Chersoni,
 543 Yu-Yin Hsu, and Simon de Deyne (eds.), *Proceedings of the Workshop on Cognitive Aspects of the
 Lexicon @ LREC-COLING 2024*, pp. 114–119, Torino, Italia, May 2024. ELRA and ICCL. URL
 544 <https://aclanthology.org/2024.cogalex-1.13/>.

- 540 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
 541 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
 542 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 543
- 544 Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing bayesian
 545 network structure learning algorithm. *Machine learning*, 65:31–78, 2006.
- 546
- 547 Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and
 548 Yun-Nung Chen. Two tales of persona in llms: A survey of role-playing and personalization. *arXiv
 549 preprint arXiv:2406.01171*, 2024.
- 550
- 551 Julius Von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel
 552 Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably
 isolates content from style. *Advances in neural information processing systems*, 34:16451–16467,
 2021.
- 553
- 554 Teng Wang, Wenhao Jiang, Zhichao Lu, Feng Zheng, Ran Cheng, Chengguo Yin, and Ping Luo.
 555 Vlmixer: Unpaired vision-language pre-training via cross-modal cutmix. In *International Conference
 556 on Machine Learning*, pp. 22680–22690. PMLR, 2022a.
- 557
- 558 Xuezhi Wang, Jason Wei, Dale Schuurmans, et al. Self-consistency improves chain of thought
 reasoning in language models. In *arXiv preprint arXiv:2203.11171*, 2023.
- 559
- 560 Y. Wang, Y. Shen, Z. Liu, P.P. Liang, A. Zadeh, and L.P. Morency. Selfattr: Self-supervised attention
 561 for multimodal attribute presentation in social media. In *Proceedings of the 30th ACM International
 562 Conference on Multimedia*, 2022b.
- 563
- 564 Yilei Wang, Jiabao Zhao, Deniz S. Ones, Liang He, and Xin Xu. Evaluating the ability of large
 565 language models to emulate personality. *Scientific Reports*, 15(1), Jan 2025. doi: 10.1038/
 s41598-024-84109-5.
- 566
- 567 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
 568 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
 569 neural information processing systems*, 35:24824–24837, 2022.
- 570
- 571 Kathryn E Wilson and Rodney K Dishman. Personality and physical activity: A systematic review
 and meta-analysis. *Personality and individual differences*, 72:230–242, 2015.
- 572
- 573 Danru Xu, Dingling Yao, Sébastien Lachapelle, Perouz Taslakian, Julius von Kügelgen, Francesco
 574 Locatello, and Sara Magliacane. A sparsity principle for partially observable causal representation
 575 learning. *arXiv preprint arXiv:2403.08335*, 2024.
- 576
- 577 H. Yang, Z. Zhang, and L. Yin. Persemon: A deep network for joint analysis of apparent personality,
 emotion and their relationship. *IEEE Transactions on Affective Computing*, 2020.
- 578
- 579 Dingling Yao, Danru Xu, Sébastien Lachapelle, Sara Magliacane, Perouz Taslakian, Georg Martius,
 580 Julius von Kügelgen, and Francesco Locatello. Multi-view causal representation learning with
 581 partial observability. *arXiv preprint arXiv:2311.04056*, 2023.
- 582
- 583 Wu Youyou, Michal Kosinski, and David Stillwell. Computer-based personality judgments are more
 584 accurate than those made by humans. *Proceedings of the National Academy of Sciences*, 112(4):
 585 1036–1040, 2015. doi: 10.1073/pnas.1418680112. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1418680112>.
- 586
- 587 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks.
 In *International Conference on Machine Learning*, pp. 7154–7163. PMLR, 2019.
- 588
- 589 Yue Yu, Tian Gao, Naiyu Yin, and Qiang Ji. Dags with no curl: An efficient dag structure learning
 590 approach. In *International Conference on Machine Learning*, pp. 12156–12166. PMLR, 2021.
- 591
- 592 Changhe Yuan, Brandon Malone, and Xiaojian Wu. Learning optimal bayesian networks using a*
 593 search. In *Twenty-second international joint conference on artificial intelligence*, 2011.
- Leslie Zebowitz. *Reading faces: Window to the soul?* Routledge, 2018.

- 594 Chao Zhang, Zichao Yang, Xiaodong He, and Li Deng. Multimodal intelligence: Representa-
 595 tion learning, information fusion, and applications. *IEEE Journal of Selected Topics in Signal*
 596 *Processing*, 14(3):478–493, 2020.
- 597
- 598 Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional
 599 independence test and application in causal discovery. *arXiv preprint arXiv:1202.3775*, 2012.
- 600
- 601 Kun Zhang, Shaoan Xie, Ignavier Ng, and Yujia Zheng. Causal representation learning from multiple
 602 distributions: A general setting. *arXiv preprint arXiv:2402.05052*, 2024.
- 603
- 604 Lecheng Zheng, Zhengzhang Chen, Jingrui He, and Haifeng Chen. Multi-modal causal structure
 605 learning and root cause analysis. *arXiv preprint arXiv:2402.02357*, 2024.
- 606
- 607 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous
 608 optimization for structure learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- 609
- 610 Yujia Zheng, Ignavier Ng, and Kun Zhang. On the identifiability of nonlinear ica: Sparsity and
 611 beyond. *Advances in neural information processing systems*, 35:16411–16422, 2022.
- 612
- 613 Klea Ziu, Slavomír Hanzely, Loka Li, Kun Zhang, Martin Takáč, and Dmitry Kamzolov.
 614 ψ dag: Projected stochastic approximation iteration for dag structure learning. *arXiv preprint*
 615 *arXiv:2410.23862*, 2024.
- 616
- 617
- 618
- 619
- 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	<i>Appendix for</i>	
649		
650	“Persona^X: Multimodal Datasets with LLM-Inferred Behavior Traits”	
651	Table of Contents:	
652		
653		
654		
655	A1 Broader Impacts and Limitations	14
656		
657	A2 Related Work	14
658		
659	A2.1 Human Behavior Trait Analysis	14
660	A2.2 Causal Discovery and Causal Representation Learning	15
661	A2.3 Multimodality and Representation Learning	16
662	A2.4 LLM Reasoning and Inference	17
663		
664		
665	A3 Details about AthlePersona and CelebPersona Datasets	18
666		
667	A3.1 Full Feature Lists	18
668	A3.2 Distribution Plots	18
669	A3.3 Missing Values	19
670		
671		
672	A4 Details about How to Select LLMs and Design Prompts	20
673		
674	A4.1 Details about How to Select LLMs (in Table ??)	20
675	A4.2 Details about the Impact of Scoring Scale (in Figure ??)	22
676		
677	A5 Details about Independent Test Results	27
678		
679	A5.1 Details about AthlePersona Dataset	28
680	A5.2 Details about CelebPersona Dataset	29
681		
682	A6 Causal Formulation, Theorems, and Proofs	29
683		
684	A6.1 Causal Model Formulation	30
685	A6.2 Proof of Theorem 3.1	30
686	A6.3 Proof of Theorem 3.2	33
687	A6.4 Proof of Theorem 3.3	34
688		
689		
690	A7 Details about Network Training for Causal Representation Learning	36
691		
692	A8 Details about Synthetic Experiments on Variant MNIST	38
693		
694	A8.1 Details about Experimental Setup	38
695	A8.2 Details about Results and Analysis	38
696		
697	A9 Details about Real-world Personality Analysis on Persona^X	38
698		
699	A9.1 Details about Experimental Setup	38
700	A9.2 Details about Results and Analysis	39
701		

702 **A1 BROADER IMPACTS AND LIMITATIONS**

703

704 Understanding human behavioral traits has broad implications for psychology, human-computer
 705 interaction, and AI personalization. This work contributes by introducing two multimodal, publicly
 706 accessible datasets (*CelebPersona* and *AthlePersona*) together with a two-layer causal anal-
 707 ysis framework, enabling the study of behavioral traits in relation to facial and biographical features.
 708 These resources provide new opportunities for exploring how observable attributes connect with
 709 higher-level human characteristics.

710 At the same time, several ethical considerations and limitations must be acknowledged. First,
 711 potential misuse is a risk: inferring behavioral or personality-related traits from public data could
 712 reinforce societal biases or be applied in harmful contexts (e.g., hiring, surveillance). To mitigate
 713 this, all data is ethically sourced, anonymized, and released only as transformed embeddings without
 714 names or raw content, ensuring consent-based, non-commercial academic use. Second, our datasets
 715 are not demographically balanced. *AthlePersona* currently includes only male athletes, while
 716 *CelebPersona* is limited to wealthy, high-visibility individuals. Findings should therefore be
 717 interpreted as population-specific rather than universal. Nonetheless, their complementary domains
 718 provide a valuable starting point for uncovering invariant causal patterns that extend across groups.

719 Overall, this work is intended to advance understanding rather than enable deterministic individual
 720 predictions. We encourage users to prioritize fairness, transparency, and privacy, and emphasize that
 721 responsible research requires ongoing dialogue about the ethical and social impacts of such datasets.
 722 Future work will focus on extending these datasets to more diverse and inclusive populations, helping
 723 reduce bias and improve generalizability.

724

725 **A2 RELATED WORK**

726

727

728 **A2.1 HUMAN BEHAVIOR TRAIT ANALYSIS**

729

730 Human behavior trait has long been central to understanding individual behavior, shaping how people
 731 reflect on their strengths, limitations, and interpersonal tendencies ?????. Traditionally, personality as-
 732 sessment has relied heavily on self-report instruments. One of the earliest was Cattell’s 16 Personality
 733 Factors (16PF), developed by Raymond Cattell to identify the fundamental traits underlying human
 734 personality (?). Another influential tool is the Eysenck Personality Questionnaire (EPQ), proposed
 735 by Hans and Sybil Eysenck, which focuses on dimensions such as extraversion, neuroticism, and
 736 psychotism (?). The Myers-Briggs Type Indicator (MBTI)(?), created by Katharine Cook Briggs
 737 and Isabel Briggs Myers, remains widely used in both social and professional domains, despite long-
 738 standing critiques regarding its limited psychometric reliability and empirical grounding. In contrast,
 739 the Big Five Personality Traits (?)—Openness, Conscientiousness, Extraversion, Agreeableness,
 740 and Neuroticism—have emerged as the prevailing framework in psychological research due to their
 741 strong empirical support and predictive power ?????.

741 Despite their utility, self-report assessments suffer from significant limitations, including susceptibility
 742 to biases such as social desirability and inconsistent self-perception. In response, computational
 743 personality analysis has gained traction, aiming to infer traits from observable signals rather than
 744 introspective questionnaires. This shift has been enabled by the abundance of digital behavioral data,
 745 allowing traits to be inferred from diverse modalities such as text (?), handwriting (?), speech (?),
 746 facial expressions (?), and online profiles (?). These methods offer scalable, non-intrusive, and
 747 real-time assessment capabilities, with applications in personalized interfaces, recruitment, and
 748 mental health diagnostics. However, most current approaches remain unimodal, limiting their ability
 749 to capture the complex interplay of cues that shape personality expression ?????.

750 To address this, recent studies have started investigating multimodal personality inference. PersE-
 751 moN (?) combined facial expression analysis with personality prediction but was constrained to
 752 controlled lab environments and small-scale datasets. SelfAttr (?) explored self-presentation behav-
 753 iors across text and images on social media, though it lacked psychometrically validated personality
 754 scores. Datasets such as CelebA (?) and FFHQ (?) have enabled large-scale facial attribute analysis
 755 but do not include personality annotations. FairFace (?) advances demographic fairness in facial
 756 datasets but similarly omits psychological dimensions.

More multimodal datasets include YouTube-Vlogs ?, FI-V2 ?, MuPTA ?, MDPE ?, and etc. Most of these data sets are built for personality prediction or detection tasks. Meanwhile, datasets explicitly built for personality research, such as myPersonality (?) and the OCEAN dataset (?), primarily rely on self-reported traits and textual data, lacking integration with visual or demographic features. These limitations underscore a persistent gap in large-scale, multimodal datasets that unify psychometric, visual, and biographical information. Our work addresses this gap by introducing two multimodal datasets—CelebPersona and AthlePersona—which combine facial, physical, and occupational features with Big Five personality assessments generated by multiple foundation models. This enables new directions in personality computing and facilitates more comprehensive and ethically-grounded AI systems capable of understanding human attributes across modalities.

A2.2 CAUSAL DISCOVERY AND CAUSAL REPRESENTATION LEARNING

Causal discovery ? from observational data has attracted considerable attention in recent decades. Constraint-based and score-based methods are two primary categories in causal discovery. Constraint-based methods, such as PC ? and FCI (?), leverage conditional independence tests (CIT; ?????) to estimate the graph skeleton and then determine the orientation. For score-based methods, the approach can vary based on the search strategy, which may involve greedy search, exact search, or continuous optimization. One typical score-based method with greedy search is Greedy Equivalent Search (GES) (?). The exact score-based methods are often time-consuming, such as dynamic programming (DP) ?, A* (?), and integer programming ?. NOTEARS (?) is the first work to cast the Bayesian network structure learning task into a continuous constrained optimization problem ?? with the least squares objective. Subsequent work GOLEM (?) adopts a continuous unconstrained optimization formulation with a likelihood-based objective. A line of works have extended NOTEARS to handle nonlinear cases via deep neural networks, such as DAG-GNN (?) and DAG-NoCurl (?). Some methods are developed to improve the computational efficiency, e.g., ψ DAG ?. In recently years, there are activa researches on causal discovery from various data constraints, including distributed data ?, heterogeneous data ?, deterministic relations ?, latent confounder and selection bias ?, and etc.

Causal representation learning (CRL) aims to recover high-level causal variables from low-level observations, bridging machine learning and causal inference (??). It generalizes classical causal discovery ? by learning structured representations that respect causal semantics. CRL methods with identifiability guarantees typically rely on additional assumptions: (1) *Functional constraints* on the data-generating process (??); (2) *Interventional or multi-environment data* that introduce distributional shifts to expose latent structure (??); (3) *Multimodal or multiview settings*, where aligned observations across modalities help identify shared causal factors through sample-level invariance (??). Recent studies unify these approaches under general invariance principles, showing that many can be seen as special cases of a broader framework (?).

Parallel to these theoretical advances, some works focus on practical CRL applications without strict identifiability. Examples include variational methods for biomedical data (?), contrastive learning for multimodal causal analysis (?), and causal discovery on neuroimaging datasets (?). In contrast, our work aims to combine both theory and application: we derive formal identifiability conditions under multi-modality multi-measurement settings and integrate these insights into a practical estimation framework for human trait analysis.

A2.3 MULTIMODALITY AND REPRESENTATION LEARNING

In the context of broader machine learning, multimodal representation learning focuses on integrating information from multiple modalities, such as text, images, and audio, to learn unified representations for downstream tasks (??). Among these methods, contrastive learning has emerged as a powerful approach, particularly for weakly supervised settings, due to its scalability and effectiveness (??????). A prominent example is CLIP model (?), which aligns text and image embeddings through contrastive objectives ???.

Unlike these methods that primarily aim for discriminative or generative performance, our work is centered on uncovering the underlying causal structure shared across modalities, with the specific goal of generating insights into personalities through principled causal representations.

810
Prompt 1. (Final Prompt for Behavior Trait Inference in AthlePersona and CelebPersona)

811
Task Description

812 Analyze the Big Five behavior traits of the individual described below. Base the analysis on publicly available information, such as direct
 813 quotes from interviews, observed public behavior, documented career patterns, and biographical details. Avoid speculation or information
 814 from unreliable gossip sources. The analysis should reflect the public persona, not a definitive psychological diagnosis or clinical evaluation.

815
Individual Information

- 816 • Name: {name}
- 817 • Gender: {gender}
- 818 • Description: {league} player, from {country} (AthlePersona) — {occupation}, from {country} (CelebPersona)

819
Instructions

- 820 1. For each of the five Big Five behavior traits (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism):
 - 821 • **Analysis:** Provide a concise (1–2 sentences) analysis. If there is sufficient public information, identify specific examples
 822 of behaviors, statements, or patterns and explain how they relate to the definition of the trait. If there is insufficient
 823 information, state that clearly.
 - 824 • **Score:** Assign a score from 0 to 3 based on the scale below.
 - 825 • **Justification:** Provide a brief (1 sentence) justification for the score, directly referencing the evidence mentioned in the
 826 analysis or the lack thereof.
- 827 2. **Summary:** After analyzing all five traits, provide a summary string containing the five scores separated by hyphens.
- 828 3. **Anonymity:** Do not explicitly mention the name of the individual in the output, use pronouns {He/His or She/Her} instead.
- 829 4. **Distinguishing Scores:** When analyzing each trait, carefully consider whether the information is insufficient (Score 0) or if it's
 830 present but indecisive (Score 2). If the available information is too sparse or vague to form any meaningful analysis, assign Score
 831 0. If there is sufficient information but it leads to an indecisive conclusion, assign Score 2.
- 832 5. **Strict Formatting:** Adhere EXACTLY to the "Expected Output Format" template below, including line breaks. Do not add any
 833 introductory or concluding remarks outside this structure.

834
Scoring Scale

- 835 • **0 = Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is
 836 unknown or unclear due to lack of data.
- 837 • **1 = Disagree** – Clear evidences contradict the trait
- 838 • **2 = Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or
 839 contradict the trait.
- 840 • **3 = Agree** – Clear evidences support the trait

841
Expected Output Format

842 **Openness:**

- 843 - Analysis: [Analysis]
- 844 - Score: [0–3]
- 845 - Justification: [Justification]

846 **Conscientiousness:**

- 847 - Analysis: [Analysis]
- 848 - Score: [0–3]
- 849 - Justification: [Justification]

850 **Extraversion:**

- 851 - Analysis: [Analysis]
- 852 - Score: [0–3]
- 853 - Justification: [Justification]

854 **Agreeableness:**

- 855 - Analysis: [Analysis]
- 856 - Score: [0–3]
- 857 - Justification: [Justification]

858 **Neuroticism:**

- 859 - Analysis: [Analysis]
- 860 - Score: [0–3]
- 861 - Justification: [Justification]

862 **Summary:** [ScoreO-ScoreC-ScoreE-ScoreA-ScoreN]

863
A2.4 LLM REASONING AND INFERENCE

864 Large Language Models (LLMs) such as GPT-4 (?), PalM (?), Gemini ?, DeepSeek ?, Qwen ?,
 865 and Claude (?) have demonstrated remarkable reasoning and inference capabilities across a wide
 866 range of tasks, including arithmetic (?), commonsense reasoning (?), text editing and generation ?,
 867 code generation (?), and scientific QA (?). These models perform zero-shot or few-shot reasoning
 868 using techniques such as chain-of-thought prompting (?), self-consistency (?), active prompting ?,

Table A1: Full Table of Features and Descriptions for AthlePersona.

AthlePersona Dataset			
Feature	Type	Description	Missing Rate (%)
Id	string	Unique identifier for each athlete	0
Height	float32	Height in centimeters	0
Weight	float32	Weight in kilograms	0
Birthyear	int32	Year of birth	0
Birthmonth	int32	Month of birth	0
Birthday	int32	Day of birth	0
League	string	Name of the athlete’s league	0
Latitude	float32	Latitude of country’s central location, transformed from the nationality	0
Longitude	float32	Longitude of country’s central location	0
Chatgpt_output	string	Full trait analysis by ChatGPT encoded in embeddings	0
Gemini_output	string	Full trait analysis by Gemini encoded in embeddings	0
Llama_output	string	Full trait analysis by LLaMA encoded in embeddings	0
Chatgpt_o to Chatgpt_n	int32	Big Five scores (OCEAN) by ChatGPT	0
Gemini_o to Gemini_n	int32	Big Five scores (OCEAN) by Gemini	0
Llama_o to Llama_n	int32	Big Five scores (OCEAN) by LLaMA	0
Final_o to Final_n	int32	Final aggregate scores for Big Five traits	0
Image_1	image	First facial image embeddings of the athlete	0

confidence-based If-or-Else prompting ? and tool-augmented reasoning (?). Despite their black-box nature, LLMs have shown emergent abilities to perform structured reasoning without explicit supervision, making them powerful general-purpose inference engines.

Recent research has explored the extent to which Large Language Models (LLMs) can infer, simulate, and even express human personality traits. For example, (?) demonstrate that models such as GPT-4 can estimate Big Five personality dimensions from user-generated text with moderate accuracy, even in zero-shot settings. Similarly, (?) show that LLMs can produce consistent personality profiles when prompted, often aligning with outputs from standardized psychometric assessments. Beyond inference, other studies examine how LLMs naturally exhibit personality-like traits in their responses. (?) introduce methods to control and elicit desired personality traits in language model outputs, while (?) analyze the emergent ability of LLMs to emulate distinct personality patterns during generation. Furthermore, recent works such as (??) utilize LLMs to annotate or assess personality traits from textual data, illustrating their growing role in computational personality research.

Our work addresses these limitations by providing multimodal datasets that unite visual, physical, demographic, and personality dimensions, with multiple model-generated assessments that enable systematic evaluation.

A3 DETAILS ABOUT ATHLEPERSONA AND CELEBPERSONA DATASETS

A3.1 FULL FEATURE LISTS

The full feature tables of AthlePersona and CelebPersona are displayed in Table ?? and ???. The complete final prompt for generating personality is shown in Prompt ??, where the blue text is the highlighted information for each individual. Summary of terms of use compliance for all different sports leagues are in Table ??.

The CelebPersona dataset contains structured information about public figures, combining demographic attributes, facial characteristics, and personality assessments. Key features include basic physical attributes (height, weight), birth details (day, month, year), and location information (latitude

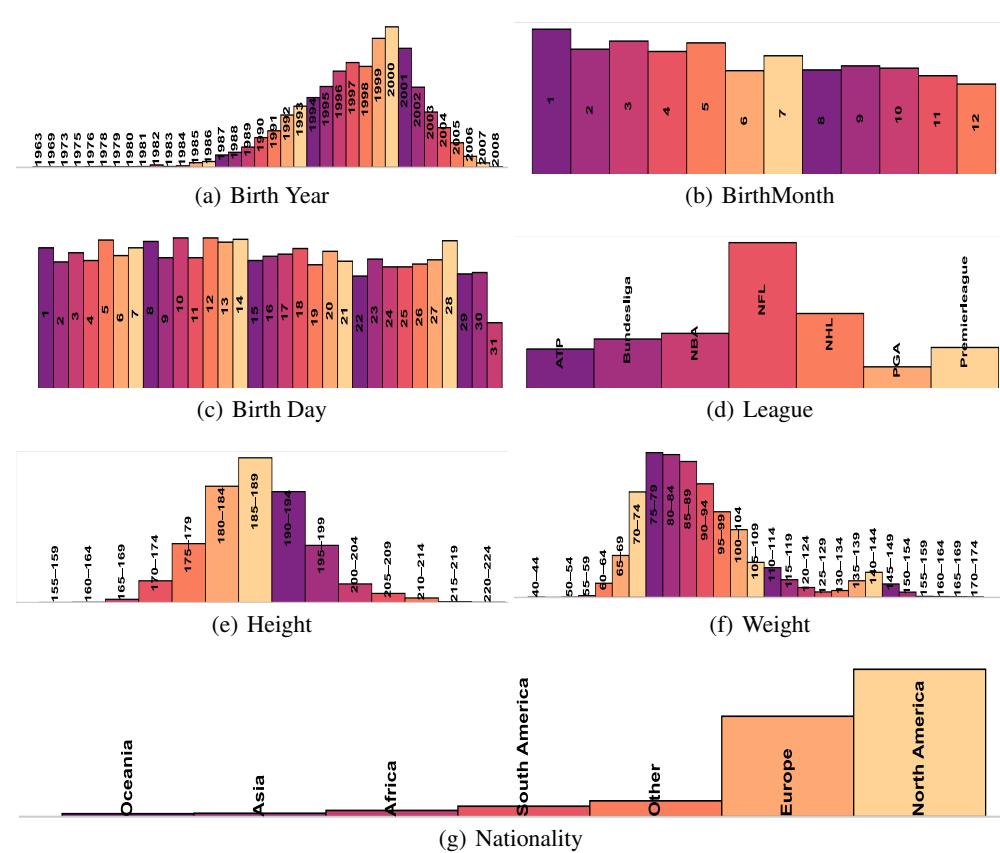


Figure A1: The distributions of the features in AthlePersona Dataset.

and longitude, transformed from their nationality). In addition, each celebrity is assigned a categorical occupation and gender label. Rich personality data is captured in the form of full-text analyses generated by ChatGPT (ChatGPT-4o-latest (2025-03-26)), Gemini (Gemini-2.5-Pro-Exp-03-25), and LLaMA (Llama-4-Maverick-03-26-Experimental), along with their respective Big Five scores (Openness, Conscientiousness, Extraversion, Agreeableness, Neuroticism, OCEAN). We choose these three models because they outperform the other models in various dimensions as shown in our preliminary experiments in Section ???. In order to represent each celebrity with one set of OCEAN personality scores, we aggregate all those three sets of 5-dimensional personality scores generated by 3 LLMs via voting, and we label those aggregated features as “Final”. Regarding the facial attributes, we manually selected 10 attributes (e.g., *Big Nose*, *High Cheekbones*) from the original 40 attributes in CelebA dataset, those selected features are most likely to present one’s inherent property in appearance, and less likely to change over short time than the others (e.g., *Heavy Makeup*, *Wearing Hat*). These binary facial attributes provide interpretable visual markers. Note that each image has a corresponding attribute value and there are multiple images per celebrity, we therefore aggregate all these attributes from different images by voting, to obtain the aggregated facial attributes. For each celebrity sample, there are at least two facial images taken from different angles, and up to 35 facial images per sample, referenced via relative file paths.

The AthlePersona dataset focuses on high-profile athletes and contains similar structure to CelebPersona, with emphasis on athletic context. It captures personal traits such as birth year, month, and day, physical measurements like height and weight, and the name of the athlete’s league. Personality descriptions and Big Five scores are again generated by ChatGPT, Gemini, and LLaMA, with final aggregated trait scores summarizing the predictions. Unlike CelebPersona, this dataset includes only a single facial image per athlete but maintains key demographic and geographic metadata. It omits facial feature annotations and categorical occupation labels, instead reflecting the athletic domain through the league information.

972 A3.2 DISTRIBUTION PLOTS
973

974 The Figure ?? shows the distribution of AthlePersona. The AthlePersona dataset is dominated by
 975 athletes born between 1985 and 2005, with a uniform spread across birth months and days. Most
 976 individuals are associated with NFL, NHL and NBA, showing a strong skew toward U.S. sports.
 977 Heights cluster around 180–199 cm, and weights around 90–109 kg, which aligns with physical
 978 norms for elite athletes in contact sports. Nationalities are overwhelmingly North American, with
 979 minimal representation from other continents, highlighting a clear Western and U.S.-centric dataset
 980 bias.

981 The Figure ?? shows the distribution of CelebPersona. The CelebPersona dataset predominantly
 982 features younger individuals, with birth years peaking between 1990–1999, and shows a balanced
 983 distribution across birth months and days. Most individuals are from North America and Europe, with
 984 underrepresentation from other continents. There is a notable occupational bias toward Entertainment,
 985 Music, and Sports, while fields like Healthcare and Academia are sparsely represented. Females
 986 slightly outnumber males. Height and weight distributions center around typical adult ranges, though
 987 outliers exist. The weight distribution peaks between 60–69 kg and 50–59 kg, with a sharp drop
 988 after 90 kg. The range 135–139 kg and higher has a minimal count. Regarding facial attributes, the
 989 majority of features are marked as absent, with a smaller subset present, particularly for traits like
 990 Oval Face and High Cheekbones. The unknown values are minimal.

991 A3.3 MISSING VALUES
992

993 As shown in the dataset features table (Table ??), the Missing Rate column indicates the proportion
 994 of unavailable or incomplete values for each feature. Despite efforts to retrieve missing infor-
 995 mation—particularly from publicly accessible sources like Wikipedia—certain attributes remain
 996 incomplete, especially those considered more private or less frequently disclosed. In the CelebPersona
 997 dataset, there are a total of 9444 data. Height and Weight have the highest number of missing entries,
 998 with 71.5% and 87% missing records respectively. Birthday and Birthmonth are missing in 2% entries
 999 each, while Birthyear is missing in 0.6% cases. Geographic coordinates (Latitude and Longitude) are
 1000 absent in 0.2% instances, and categorical attributes such as Occupation_Num and Gender_Num have
 1001 0.05% and 0.2% missing values, respectively.

1002 In contrast, the AthlePersona dataset (Table??) has been fully cleaned by removing all rows that
 1003 contain any missing values. Prior to finalization, any entry with incomplete demographic, geographic,
 1004 or profile information was excluded to ensure consistency. As a result, AthlePersona contains
 1005 no missing values, making it readily usable for downstream analysis without requiring additional
 1006 preprocessing or imputation.

1007
1008 A4 DETAILS ABOUT HOW TO SELECT LLMs AND DESIGN PROMPTS
1009

1010 Throughout this paper, we rely on large language models (LLMs) to generate human personality.
 1011 Fortunately, Benefiting from the recent explosion in the size and availability of LLMs ????, some
 1012 research has shown that personality measurements in the outputs of some LLMs under specific
 1013 prompting configurations are valid and reliable ??. We built our datasets from a number of
 1014 professional athletes and celebrities, based on the facts that they are famous and it is likely to have
 1015 sufficient information about them online.

1016 In Section ??, we present some experiments to show how to select LLMs for personality generation,
 1017 and also how to design the prompts. In the following, we will demonstrate more details.

1018
1019 A4.1 DETAILS ABOUT HOW TO SELECT LLMs (IN TABLE ??)

1020
1021 **Model Choices in Table ??.** We present the comparative evaluations on 10 of the state-of-the-art
 1022 LLMs on our two persona datasets. Initially, we choose the Top 10 models based on the Arena
 1023 leaderboard ?. To enhance diversity, we also included Qwen2.5-Max and QwQ-32B ? from Alibaba,
 1024 both noted for their strong reasoning capabilities. We list all those 12 LLMs and summarize them in
 1025 Table ??, including the model name, the company name, the arena score, the API input and output
 price per million tokens, and whether it is used by us for further analysis in Table ?. Specifically,

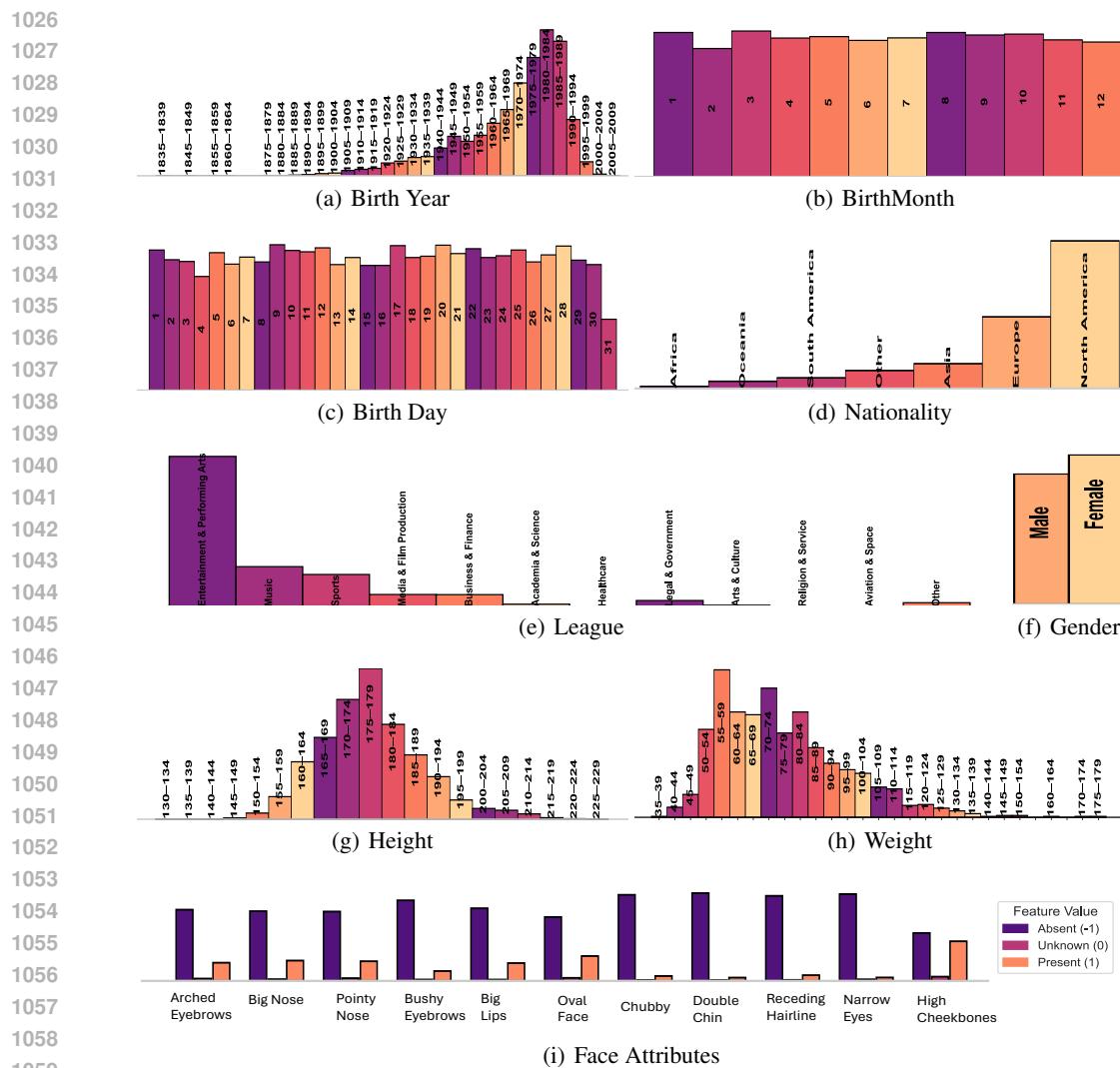


Figure A2: The distributions of the features in CelebPersona Dataset.

GPT-4.5-Preview ? was excluded due to prohibitively high API costs where the input and output API prices were \$75 and \$150 per million tokens, respectively. Gemini-2.0-Pro-Exp-02-05 ? was omitted due to inaccessibility, it was merged to the latest Gemini-2.5-Pro-Exp-03-25 model. Therefore, in the end, we only considered 10 LLMs, as shown in Table ??.

Evaluation Metrics in Table ??. We list 8 evaluation metrics in the experimental results. For each dataset, we randomly sampled 100 individuals, conducted 100 LLM queries in total, and reported the average results. The query prompt is almost the same as the Prompt ??, except that here we considered 5-level (i.e., strong disagree, disagree, neutral, agree, strongly agree) for scoring scale instead of 3-level. As for each evaluation metric, here are detailed explanations:

Generation Time (GT) measures the computational efficiency of each large language model by recording the average inference time required to produce responses. This metric is quantified in seconds and provides insight into the practical usability of different models, with lower values indicating faster processing speeds.

Missing Rate (MR) quantifies the frequency at which language models fail to provide the requested scoring output due to limitations in their knowledge base. This metric is calculated as the percentage

1080

Prompt 2. (Evaluation Prompt for Context Consistency and Factual Accuracy)

1081

Evaluation Task

1082

You are an expert evaluator for behavior trait analysis results generated by LLMs. You will analyze the following output and evaluate it on two specific criteria.

1083

1084

Input Information

1085

1086

1087

1088

1089

1090

1091

1092

Evaluation Criteria

1093

1. Context Consistency [0/1]

1094

1095

1096

1097

1098

Check if each of the Big Five trait analyses is consistent with the assigned score (0-5)

Check if the justification for each score aligns with the analysis

Score 1 if all analyses are internally consistent with their scores and justifications

Score 0 if any inconsistencies exist (e.g., describing high extraversion traits but giving a score of 2)

1099

2. Factual Accuracy Assessment [0/1]

1100

1101

1102

1103

1104

1105

1106

Required Output Format

1107

1108

1109

1110

1111

1112

of instances where the model cannot generate a proper response (with output score 0 - Unknown), highlighting gaps in the model's capability to handle certain types of queries or domains.

1113

Indecisive Rate (IR) captures the proportion of responses where models express uncertainty or provide neutral answers rather than definitive judgments (with output score 2 - Neutral). This metric reflects the model's confidence level and willingness to make clear assessments, with higher rates indicating more cautious or uncertain behavior.

1114

1115

1116

1117

1118

Privacy Preservation (PP) evaluates the model's ability to protect individual identities by effectively anonymizing personal information in its responses. This metric assesses how well the model handles sensitive data and maintains privacy standards while still providing meaningful analysis. For each response, if there contains any individual name information, return 0, otherwise return 1.

1119

1120

1121

1122

1123

1124

1125

1126

1127

Output Formatting (OF) measures adherence to specified response structure and format requirements. This metric evaluates whether the model consistently follows given instructions regarding how responses should be organized and presented, ensuring usability and consistency. For each response, if it absolutely follows the given instructions and the output template format, return 1, otherwise return 0.

1128

1129

1130

1131

Context Consistency (CC) assesses the internal coherence between different components of the model's response, specifically examining alignment between the analysis, assigned score, and provided justification.

1132

1133

Factual Accuracy (FA) measures the absence of factual errors in the model's output, evaluated through cross-validation using mutual critique between different language models. This metric is crucial for determining the reliability and trustworthiness of the generated content.

1134

Prompt 3. (Comparison on Scoring Scale for Different Prompts)

1135

1136

1137

1138

1139

1140

1141

1142

1143

[Number-L3-Inc]

- **0 = Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **1 = Disagree** – Clear evidences contradict the trait
- **2 = Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **3 = Agree** – Clear evidences support the trait

[Number-L3-Dec]

1144

1145

1146

1147

1148

1149

1150

- **0 = Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **3 = Disagree** – Clear evidences contradict the trait
- **2 = Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **1 = Agree** – Clear evidences support the trait

[Text-L3-Inc]

1151

1152

1153

1154

1155

1156

1157

- **Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **Disagree** – Clear evidences contradict the trait
- **Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **Agree** – Clear evidences support the trait

[Number-L5-Inc]

1158

1159

1160

1161

1162

1163

1164

1165

1166

- **0 = Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **1 = Strongly Disagree** – Clear evidences contradict the trait
- **2 = Disagree** – Some evidences contradict the trait
- **3 = Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **4 = Agree** – Some evidences support the trait
- **5 = Strongly Agree** – Clear, consistent evidences support the trait

[Number-L5-Dec]

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

- **0 = Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **5 = Strongly Disagree** – Clear evidences contradict the trait
- **4 = Disagree** – Some evidences contradict the trait
- **3 = Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **2 = Agree** – Some evidences support the trait
- **1 = Strongly Agree** – Clear, consistent evidences support the trait

[Text-L5-Inc]

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

- **Insufficient information** – Not enough reliable public information to assess the trait. The trait's presence or absence is unknown or unclear due to lack of data.
- **Strongly Disagree** – Clear evidences contradict the trait
- **Disagree** – Some evidences contradict the trait
- **Neutral** – Evidence is mixed or the trait is not prominent. There is enough information, but it does not strongly support or contradict the trait.
- **Agree** – Some evidences support the trait
- **Strongly Agree** – Clear, consistent evidences support the trait

Note that both CC and FA metrics were evaluated through a generator-evaluator manner by 4 different evaluator LLMs (Gemini-2.5-Pro-Exp-03-25, Llama-4-Maverick-03-26-Experimental, ChatGPT-4o-latest(2025-03-26), and Grok-3-Preview-02-24), to ensure logical consistency within responses. Basically, we collect the generated personality analysis output by 10 generator LLMs, and feed into other 4 evaluator LLMs. The evaluator LLMs will return 0 (indicating No) or 1 (indicating Yes). The evaluation prompt is presented in Prompt ???. There are mainly two reasons why we do not use human evaluators but instead choosing LLM evaluators: (1) First, human evaluator is expensive and costly; (2) Second, except loyal fans, most people may not have an in-depth understanding about a celebrity or athlete. To that end, LLMs probably have seen more information about certain celebrity or athlete than normal human in general. Therefore, for factual accuracy evaluation, it is reasonable to use LLM evaluators. Note that in this way, we aim to point out any statement which absolutely violates the factuality or commonsense. As for evaluating context consistency, it turns out to be a text interpretation task, it is also reasonable to apply LLMs.

Overall Score (OS) provides a comprehensive performance measure by calculating the average of all evaluation metrics except Generation Time. The score calculation is:

$$OS = \frac{1}{6} \times [PP + OF + CC + FA + (1 - MR) + (1 - IR)]. \quad (1)$$

It offers a holistic view of each model’s capabilities across the various assessment dimensions.

Analysis. Table ?? presents a comparative evaluation of 10 LLMs. ChatGPT-4o-Latest ? and Gemini-2.5-Pro ? achieved the highest overall scores. Performance is consistently stronger on CelebPersona than on AthlePersona, indicating that assessing athlete personalities is more challenging. This is particularly reflected in the higher MR on AthlePersona, which is possibly due to the limited public information available for younger or less prominent athletes. While GT varies substantially across models, both PP and OF are consistently strong. IR differs notably, e.g., 0.46 for Qwen-Plus ? while 0.11 for DeepSeek-R1 ?, suggesting significant variation in models’ confidence calibration.

A4.2 DETAILS ABOUT THE IMPACT OF SCORING SCALE (IN FIGURE ??)

Prior research has demonstrated the importance of prompt engineering strategies in enhancing LLM performance across various tasks ??????. These foundational studies have established that prompt structure and presentation significantly influence model outputs, particularly in psychological assessment applications. Building on this foundation, we systematically investigate how variations in *scoring scale format* affect the consistency of LLM-generated personality assessments across the Big Five traits, with implications for reliable automated psychological evaluation.

How to Design Prompts? As shown in Fig. ??, the radar and box plots in the top and middle illustrate the extent of *intra-prompt* variability across the Big Five traits, while the bottom panel reports Manhattan distances between prompts to capture *inter-prompt* differences. Across both CelebPersona and AthlePersona datasets, Llama-4-Maverick ? stands out for its highly stable outputs, followed by Gemini-2.5-Pro ?. In contrast, Qwen2.5-Max ? tends to produce the most variable results. Among these prompts, the “Number-L3-Inc” format consistently yields the lowest variance, suggesting that coarse, numerically formatted 3-point scales help LLMs produce more deterministic responses. Conversely, more complex prompts, especially those using Level-5 textual scales, lead to noticeably higher variability. Taken together, these findings suggest that prompt design, particularly scale granularity and formatting, plays a critical role in shaping the reliability of LLM-based personality assessment.

Experimental Design and Methodology in Figure ??). We evaluated five top-performing LLMs using a structured prompt format [Number/Text] { [L3/L5] { [Inc/Dec]}, where elements specify response type (numerical vs. textual), scale granularity (3-level vs. 5-level), and ordering (increasing vs. decreasing). We list all different scoring scale in different prompts in Prompt ??.

This systematic approach enables comprehensive analysis of how different formatting choices interact to influence model behavior. Each model was tested across 100 trials per prompt format on both CelebPersona and AthlePersona datasets, with temperature set to 0 to reduce stochastic variability (even though temperature 0 will still have output variation) and isolate prompt-related effects.

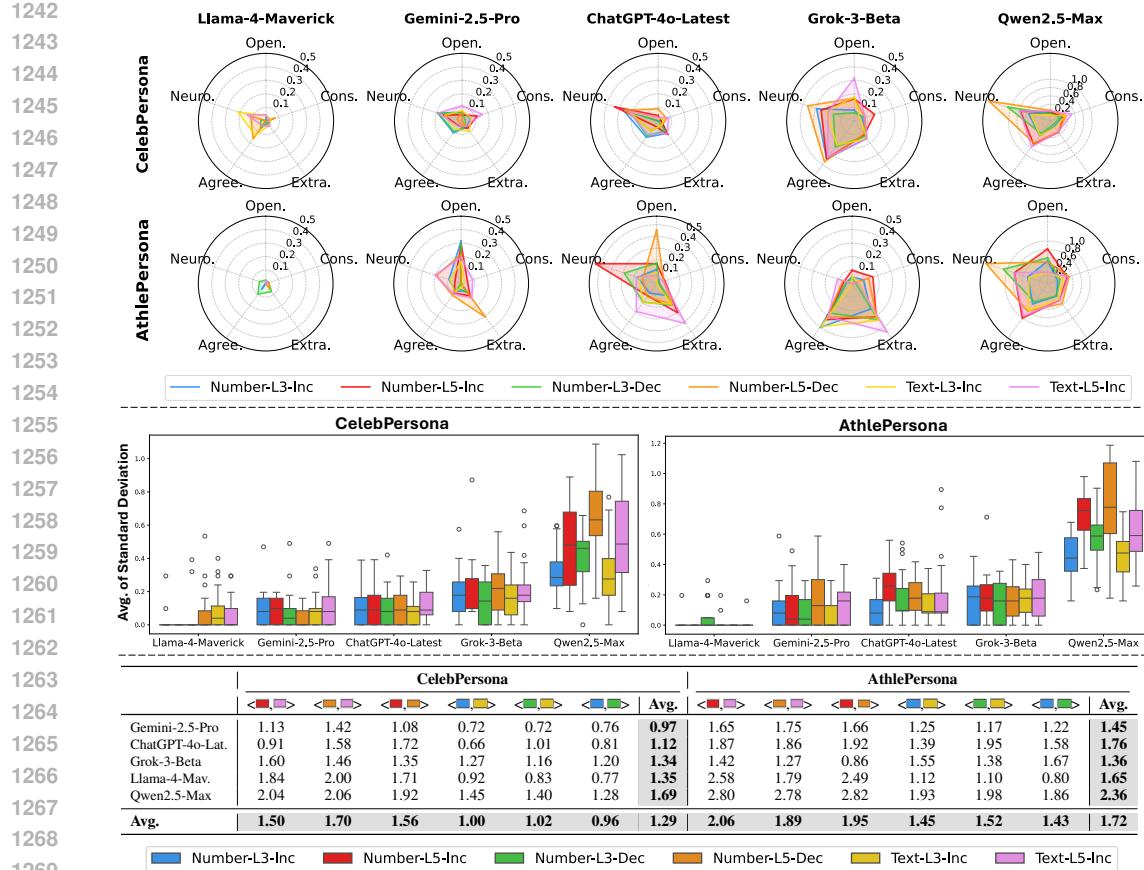


Figure A3: **Impact of scoring scale on LLM consistency.** **Top:** Radar plots show the standard deviation (std) of Big Five trait scores across repeated runs under different prompt formats, separately for each model (by column) and dataset (by row). **Middle:** Box plots summarize the average of std across Big Five personality traits, highlighting *intra-prompt* variability. **Bottom:** Manhattan distances between two prompt pairs quantify *inter-prompt* variability. Refer to § ?? for more setup and result analysis.

Consistency was quantified using standard deviation (std) of personality scores across repeated runs, providing direct measures of output stability.

Comprehensive Analysis Framework. Our analysis encompasses three complementary perspectives as shown in Figure ???: (1) Top: trait-specific variability patterns through radar plots, (2) Middle: aggregate consistency measures via box plot distributions, and (3) Bottom: inter-prompt relationship quantification using Manhattan distance matrices. This multi-faceted approach provides both granular insights into individual trait reliability and broader patterns in prompt format effectiveness.

Model Performance Hierarchy and Stability Patterns. The analysis reveals a clear performance hierarchy among evaluated models. Llama-4-Maverick demonstrates exceptional consistency with standard deviations consistently below 0.2 across all prompt formats and personality traits, forming tight, regular polygons in radar plots that indicate robust internal mechanisms for maintaining consistent assessments. The model's box plots show minimal variability between prompt formats with few outliers, suggesting sophisticated handling of diverse input structures.

Gemini-2.5-Pro occupies an intermediate position with generally low variability but occasional sensitivity to specific prompt formats, evidenced by longer box plot whiskers and more distributed quartiles. The model shows particular stability with numerical scales while demonstrating increased variance with textual scales, indicating format-dependent reliability patterns. ChatGPT-4o-Latest exhibits moderate consistency overall but with notable prompt-dependent variations, particularly visible through outliers in box plot distributions. While generally reliable, certain prompt-model-trait combinations produce unexpectedly high variability, suggesting sensitivity to specific formatting

1296 choices. Grok-3-Beta shows concerning instability, particularly in AthlePersona where some prompt
 1297 formats yield standard deviations exceeding 0.8. Wide interquartile ranges indicate dramatic con-
 1298 sistency variations depending on prompt format, with pronounced radar plot irregularities revealing
 1299 trait-specific vulnerabilities. Qwen2.5-Max consistently ranks as the least reliable model, exhib-
 1300 iting high median standard deviations and extensive outliers reaching above 1.0. The model’s radar
 1301 plots often show expanded, irregular shapes indicating inconsistent performance across traits, with
 1302 Manhattan distances exceeding 2.0 for complex formats.

1303 **Trait-Specific Consistency Patterns.** The radar plot analysis reveals compelling trait-specific
 1304 reliability patterns. Openness emerges as the most stable trait across nearly all models and prompt
 1305 formats, consistently showing standard deviations below 0.3. This stability suggests that LLMs
 1306 demonstrate inherent consistency when evaluating creative and intellectual characteristics, possibly
 1307 due to clearer linguistic markers for openness-related traits in training data.

1308 Neuroticism presents notable dataset dependency, showing moderate stability in CelebPersona but
 1309 considerably higher variability in AthlePersona, particularly for less stable models where standard de-
 1310 viations can exceed 1.0. This context-dependent pattern indicates that evaluation domain significantly
 1311 influences how models interpret emotional stability markers. Extraversion and Agreeableness exhibit
 1312 intermediate variability levels with distinct model-specific patterns. The geometric shapes formed by
 1313 different prompt formats in radar plots reveal systematic differences: simpler formats tend to create
 1314 smaller, more regular polygons, while complex textual formats often produce irregular, expanded
 1315 shapes indicating inconsistent cross-trait performance.

1316 **Format Optimization and Complexity Trade-offs.** The Number-L3-Inc format consistently yields
 1317 the lowest variance across models and datasets, demonstrating that simple numerical 3-level scales en-
 1318 hance deterministic responses. Box plot analyses show this format produces the tightest distributions
 1319 with minimal outliers across all models. Manhattan distance matrices reveal that Number-L3-Inc and
 1320 Number-L3-Dec formats show consistently low inter-prompt distances (often below 1.0), indicating
 1321 that scale direction has minimal impact when using simple numerical formats.

1322 Conversely, textual 5-level formats (Text-L5-Inc/Dec) produce significantly higher variability, with
 1323 standard deviations often exceeding 0.5 and Manhattan distances reaching above 2.0 between prompt
 1324 pairs. This indicates that textual formats not only increase intra-prompt variability but fundamen-
 1325 tally alter response distributions compared to numerical approaches. The increased granularity of
 1326 5-level scales appears to introduce additional decision boundaries that models interpret inconsistently.
 1327 Number-L5 formats show intermediate complexity, exhibiting distances that fall between L3 numeri-
 1328 cal formats and textual formats. This suggests that 5-level scales represent a transitional complexity
 1329 level—more challenging than 3-level scales but not as fundamentally different as textual ones.

1330 **Cross-Dataset Insights and Domain Effects.** Systematic comparison between CelebPersona and
 1331 AthlePersona reveals important domain-dependent patterns. AthlePersona generally produces higher
 1332 standard deviations and inter-prompt distances across most models, suggesting that athlete personality
 1333 assessment presents inherent challenges for LLMs. This pattern may reflect training data biases,
 1334 where celebrity personalities are more extensively documented in text corpora compared to athlete
 1335 psychological profiles, leading to less robust assessment capabilities in athletic contexts.

1336 **Implications and Final Model Selection.** *These findings challenge conventional assumptions
 1337 about measurement precision in automated assessment contexts. Counter-intuitively, reducing scale
 1338 granularity and employing numerical rather than textual formats substantially improves reliability,
 1339 suggesting that cognitive complexity reduction outweighs precision benefits of more detailed scales.
 1340 Based on our comprehensive analysis across multiple evaluation dimensions, we made the following
 1341 strategic selections for our personality generation framework: After careful consideration of the
 1342 consistency patterns, trait-specific reliability, and cross-dataset performance, we chose the **Number-
 1343 L3-Inc format** as our standardized prompt structure. This format demonstrated the lowest variance
 1344 across all models and datasets, with standard deviations consistently below 0.3 and minimal inter-
 1345 prompt distances, ensuring maximum reliability in automated personality assessment.*

1346 *For model selection, we adopted a multi-model approach incorporating **Llama-4-Maverick**, **ChatGPT-
 1347 4o-Latest**, and **Gemini-2.5-Pro**. Llama-4-Maverick serves as our primary model due to its exceptional
 1348 consistency ($std \approx 0.2$) across all traits and formats. ChatGPT-4o-Latest provides complementary
 1349 reliability with moderate consistency and broad accessibility, while Gemini-2.5-Pro offers additional
 validation particularly for numerical format processing. This ensemble approach leverages the*

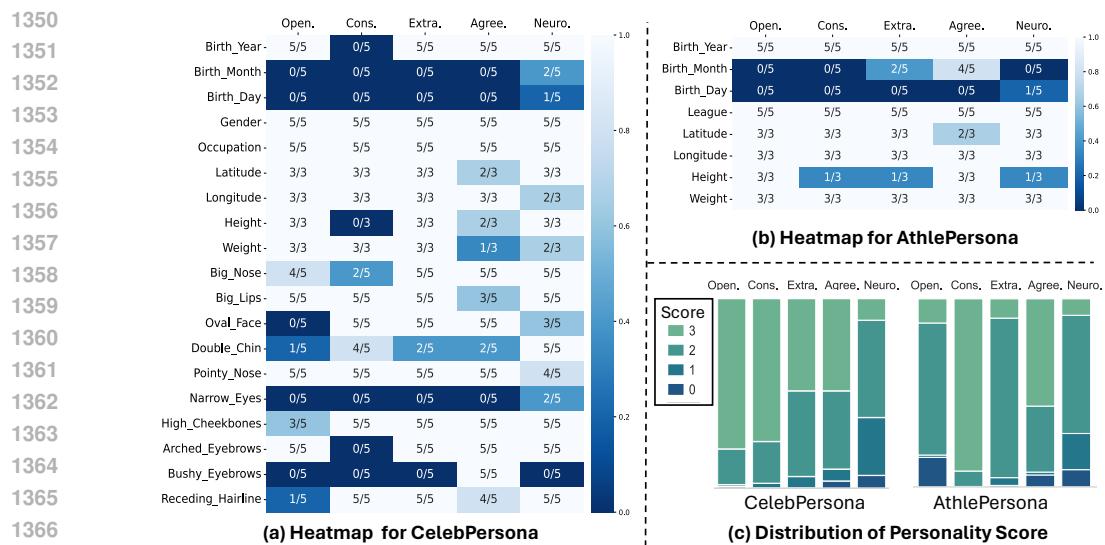


Figure A4: **Independence test (IT) heatmap and personality score distribution.** (a)/(b) present heatmaps of significant IT results between Big Five personality and other features for CelebPersona and AthlePersona, respectively. Each cell shows “ x/y ”, where x is the number of significant test methods with $p < 0.05$, and y is the total number of test methods applied. Lighter areas indicate stronger dependence. (c) shows the distribution of Big Five personality scores across both datasets.

strengths of multiple models while mitigating individual model limitations observed in our analysis. Notably, we excluded Grok-3-Beta and Qwen2.5-Max from our final selection due to their concerning instability patterns, with standard deviations frequently exceeding 0.8 and inconsistent cross-trait performance that could compromise assessment reliability.

The observed trait-specific and dataset-dependent variations underscore the critical importance of careful prompt design in LLM-based psychological evaluation systems. The convergent evidence across radar plots, box plot distributions, and distance matrices demonstrates that prompt engineering represents a fundamental factor in determining assessment reliability, with implications extending beyond personality evaluation to broader automated psychological assessment applications.

A5 DETAILS ABOUT INDEPENDENT TEST RESULTS

To evaluate independence relationships across different variable types in our analysis, we employ five statistical testing methods. For discrete variables, we utilize two classical approaches: the Chi-square test ?, which evaluates statistical independence between categorical variables, and the G-square test ?, a likelihood-ratio variant that demonstrates improved robustness in small sample scenarios. For continuous and mixed variable types, we implement three kernel-based methods: the Hilbert-Schmidt Independence Criterion (HSIC) ?, which measures dependence in high-dimensional data using reproducing kernel Hilbert spaces; the Randomized Conditional Independence Test (RCIT) ?, a non-parametric approach that employs randomized Fourier features to approximate kernel-based conditional independence testing; and the Kernel-based Conditional Independence Test (KCI) ?, which extends HSIC methodology for testing conditional independence in complex data structures. This comprehensive suite of methods enables robust independence testing across diverse data types encountered in our experimental framework. A dependency is deemed significant if $p < 0.05$, and each cell in Fig. ??(a)/(b) shows the number of methods that detect such significance, and we summarize these 5 methods in Table ?? and Table ??.

A5.1 DETAILS ABOUT ATHLEPERSONA DATASET

Figure ?? presents heatmaps of p-values from different statistical independence. The Chi-Square Test (CSQ) and G-Square Test (GSQ) show remarkably similar patterns, which is expected given their shared theoretical foundation for categorical variables. Overall, the independence test analysis reveals

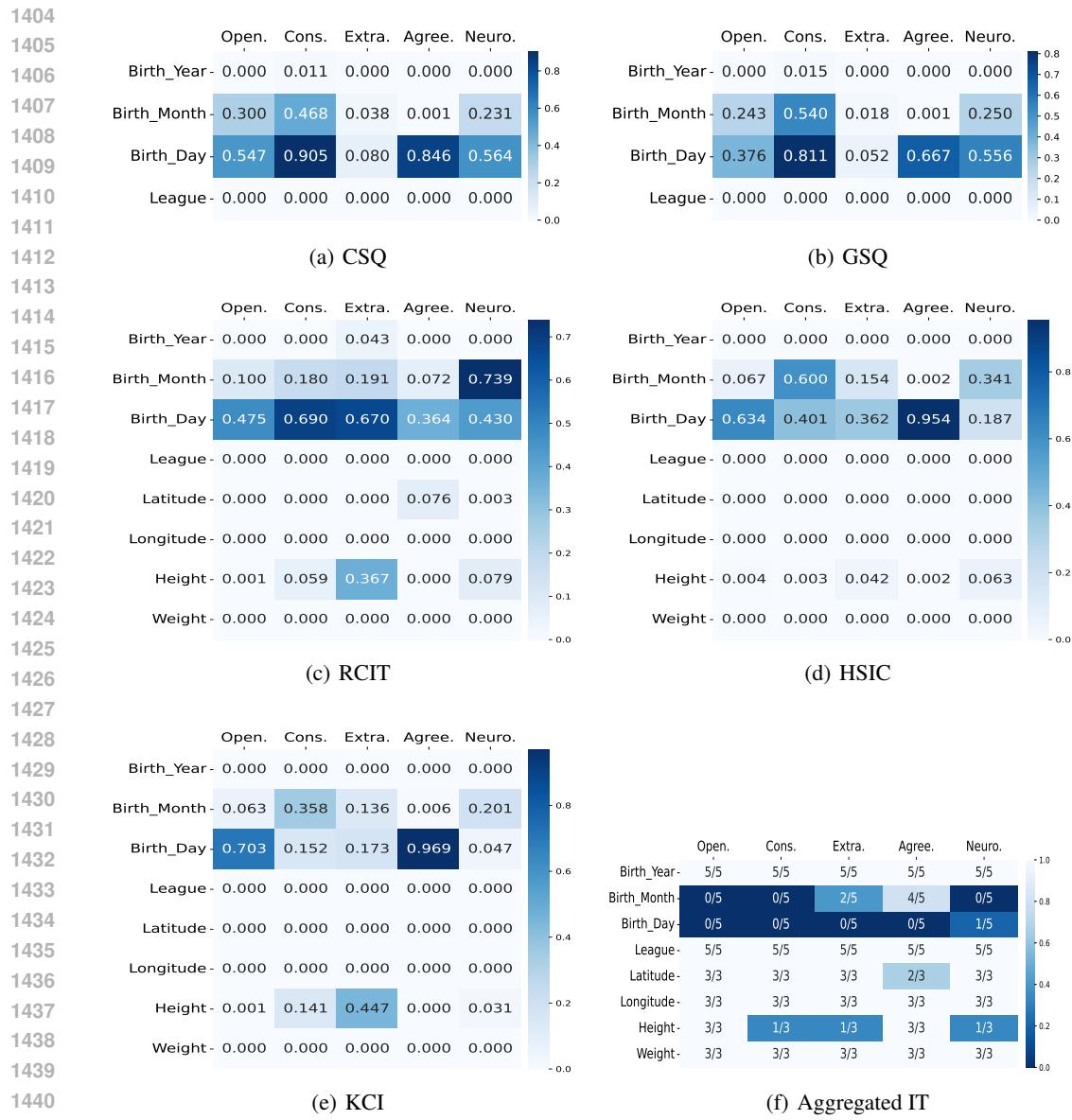


Figure A5: AthlePersona: Heatmap of P-value obtained from different independence test. (a) Chi-Square Test (b) G-Square Test (c) Randomized Conditional Independent Test (d) Hilbert-Schmidt Independence Test (e) Kernel-based Conditional Independence Test (f) Aggregated results over all above 5 methods.

limited but significant demographic-personality dependencies in the AthlePersona dataset. Most relationships show p-values well above the 0.05 significance threshold, indicating statistical independence between demographic features and personality traits. However, notable exceptions include birth year and league's strong dependence with all big five personality traits, birth month associations with Agreeableness in the CSQ test ($p = 0.001$), which represents the strongest dependency detected. Birth day shows somewhat clear independence with personality in most methods. The kernel-based methods (RCIT, HSIC, KCI) generally produce lower p-values, indicating stronger evidence for dependencies. Most relationships show p-values between 0-0.01, suggesting statistical dependence among variables such as birth year, league, latitude, longitude, weight, and the Big Five personality traits. Interestingly, weight is more dependent on openness, agreeableness, and neuroticism, while being more independent of conscientiousness and extraversion.

Trait-specific analysis reveals that most Big Five personality dimensions operate independently of the measured demographic factors in athletic populations. Agreeableness shows the most consistent evidence of demographic sensitivity, particularly with birth timing variables, though significant relationships ($p < 0.05$) remain infrequent across methods. Openness, Conscientiousness, Extraversion, and Neuroticism demonstrate predominantly dependent relationships with demographic features, with p-values typically smaller than 0.05 across most variable-method combinations. Particularly, league affiliation, geographic coordinates (latitude, longitude), and birth year show consistent results, with most methods yielding very low p-values (near 0.000) suggesting dependence, while birth month and birth day produce high p-values indicating independence.

The multi-method validation approach reveals important methodological insights about dependency detection reliability. Classical categorical tests (CSQ, GSQ) occasionally detect marginal associations that kernel-based methods (RCIT, HSIC, KCI) fail to identify, suggesting method-specific sensitivities rather than robust dependencies. The independence test heatmap shows mixed results: some variables like birth month, birth day, and height demonstrate low consensus scores (0-2 out of 5 methods achieving $p < 0.05$), indicating weak or inconsistent dependencies. However, several variable-trait combinations achieve moderate to high consensus scores, primarily involving league, latitude, longitude, and weight. This pattern suggests a nuanced relationship where certain demographic factors (geographic and league-related variables) show more consistent associations with personality traits in athletic populations than temporal or physical characteristics.

The dependencies between Big Five personality traits and league, latitude, longitude, and weight in athletic populations likely reflect a complex interplay of self-selection, environmental influences, and sport-specific demands. League affiliations may attract distinct personality profiles—team sports favoring extraversion and agreeableness for collaboration, while individual sports might select for conscientiousness and controlled neuroticism. Geographic variables (latitude/longitude) capture regional cultural differences in values like individualism versus collectivism, as well as environmental factors such as climate that research has linked to personality development. Weight dependencies may emerge through multiple pathways: conscientiousness influencing self-regulation of diet and exercise, neuroticism affecting stress-related eating behaviors, openness driving willingness to experiment with training regimens, and sport-specific body type requirements that indirectly link physical characteristics to the personality traits favored in those sports. These relationships represent genuine demographic-personality associations rather than statistical noise because they align with theoretically plausible mechanisms involving cultural adaptation, environmental pressures, and the mutual influence between personality traits and lifestyle choices in elite athletic contexts.

A5.2 DETAILS ABOUT CELEBPERSONA DATASET

Figure ?? shows heatmaps of p-values from different statistical independence tests evaluating the relationship between facial/demographic features and Big Five personality traits in the CelebPersona dataset. Features like birth year, gender, occupation, latitude, longitude, pointy nose and big lips frequently show strong associations with Big Five traits. In contrast, attributes like birth day, narrow eyes and bushy eyebrows generally appear independent of personality.

The CelebPersona dataset reveals several robust dependency patterns with p-values consistently below 0.05 across multiple methods. Birth timing variables demonstrate the strongest dependencies: birth day shows significant associations with openness, conscientiousness, and extraversion across kernel-based methods, suggesting developmental timing effects on personality formation. Birth month exhibits dependencies with conscientiousness and moderate associations across other traits. Among facial features, big nose demonstrates consistent dependencies with conscientiousness across kernel methods, while bushy eyebrows shows significant associations with openness and extraversion. Weight exhibits notable dependencies with agreeableness and neuroticism, indicating body composition-personality linkages. Narrow eyes shows dependencies with conscientiousness and agreeableness, while oval face demonstrates associations with neuroticism and other traits.

The aggregated IT results confirm these dependencies with higher consensus scores for birth day (3-4/5 methods), bushy eyebrows (3-4/5 methods), and weight (2-3/5 methods), indicating genuine associations rather than statistical noise. Classical methods (CSQ, GSQ) detect fewer significant relationships, suggesting that non-linear dependency structures dominate celebrity personality-morphology associations. These findings support evolutionary psychology theories linking facial

morphology to personality traits, particularly the relationship between eyebrow prominence and openness/extraversion, and nose characteristics with conscientiousness. The effects of the timing of the birth may reflect seasonal developmental influences or cohort effects specific to the career trajectories of the entertainment industry, where certain combinations of personality and timing of the birth provide advantages in celebrity achievement.

The dependency patterns in celebrity populations reveal intriguing domain-specific insights that distinguish them from general populations. The pronounced birth timing effects, particularly the strong associations between birth day and multiple personality traits, suggest that developmental timing may interact with entertainment industry selection pressures in unique ways. Celebrities born on certain days may possess personality configurations that enhance their ability to navigate public scrutiny, media attention, and performance demands. The facial feature dependencies present a complex picture of appearance-personality relationships: the consistent association between bushy eyebrows and openness/extraversion aligns with research on facial masculinity and dominance signaling, while the nose-conscientiousness relationship may reflect underlying genetic correlations between facial development and self-regulatory capacity. Weight dependencies with agreeableness and neuroticism indicate that body image management, a critical aspect of celebrity careers, may both influence and be influenced by personality traits related to social harmony and emotional stability. The higher dependency rates detected by kernel methods compared to classical approaches suggest that celebrity personality-morphology relationships involve complex, non-linear interactions that traditional statistical methods fail to capture, possibly reflecting the multifaceted nature of public persona where appearance, personality, and career success form intricate feedback loops.

A6 CAUSAL FORMULATION, THEOREMS, AND PROOFS

In this section, we will show the causal model formulation firstly, then we will present all the theorems, and their proofs. In Theorem 1, We begin by showing how the modality-specific latent subspaces $[\mathbf{z}_m, \mathbf{s}]$, where \mathbf{z}_m is modality-specific latent variables and \mathbf{s} is modality-shared latent variables, can be recovered in a nonparametric manner using multiple measurements. Building on this result, then in Theorem 2, we demonstrate the identifiability of the shared latent variable \mathbf{s} by leveraging the information across multiple modalities. Finally in Theorem 3,, conditioned on the recovered \mathbf{s} , we establish the identifiability of each modality-specific latent variable \mathbf{z}_m up to minor indeterminacies, i.e., component-wise identifiability with an inner-modality permutation. The logical dependencies among the theorems are summarized in the flowchart as shown in Figure ??.

A6.1 CAUSAL MODEL FORMULATION

Data-generating processes. Let $\mathbf{x} := [\mathbf{x}_1, \dots, \mathbf{x}_M]$ be a set of observations/measurements from M modalities, where $\mathbf{x}_m \in \mathbb{R}^{d_m}$ represents the observation from modality m with dimensionality d_m . Let $\mathbf{z} = [\mathbf{z}_1, \dots, \mathbf{z}_M]$ be the set of causally related latent variables underlying m -th modalities. Specifically, the data generation process (Figure ??) can be formulated as

$$z_{m,i} := g_{z_{m,i}}(\text{Pa}(z_{m,i}), \mathbf{s}, \epsilon_{m,i}), \quad (\text{latent causal relations}) \quad (2)$$

$$\mathbf{x}_m := g_{\mathbf{x}_m}(\mathbf{z}_m, \boldsymbol{\eta}_m), \quad (\text{generating functions}) \quad (3)$$

where we denote the parents of a variable with $\text{Pa}(\cdot)$. Since we allow for general causal relations within each modality and across multiple modalities, $\text{Pa}(\cdot)$ potentially returns latent variables across multiple modalities. Additionally, we allow the shared latent variable \mathbf{s} generally governing the modality-specific latent variables \mathbf{z}_m . The differentiable function $g_{\mathbf{z}}$ encodes the latent causal graph connecting latent components and its Jacobian matrix $\mathbf{J}_{g_{\mathbf{z}}}$ can be permuted into a strictly triangular matrix. We use $\epsilon_{m,i}$ to denote the exogenous variable for $z_{m,i}$ and exogenous variables are mutually independent. We use $\boldsymbol{\eta}_m$ to denote modality-specific information independent of other components.

Identifiability Definition. As mentioned previously, our aim was to learn the latent variables underlying each modality and their causal relations. Formally, for two specifications $\theta := \{g_{\mathbf{x}_m}, g_{\mathbf{z}_m}, p(\mathbf{s}), p(\epsilon_m), p(\boldsymbol{\eta}_m)\}_{m=1}^M$ and $\hat{\theta} := \{\hat{g}_{\mathbf{x}_m}, \hat{g}_{\mathbf{z}_m}, \hat{p}(\mathbf{s}), \hat{p}(\epsilon_m), \hat{p}(\boldsymbol{\eta}_m)\}_{m=1}^M$ of the data-generating process Eq. equation ?? and Eq. equation ?? that fit the marginal distribution $p(\mathbf{x})$, we would like to show that: given the same \mathbf{x} value, each latent component $\hat{z}_{m,i}$ is equivalent to its

1566 counterpart $z_{m,i}$ up to an invertible map $h_{m,i}$, i.e., $\hat{z}_{m,i} = h_{m,i}(z_{m,i})$. This property is known as
 1567 identifiability.
 1568
 1569
 1570

A6.2 PROOF OF THEOREM 3.1

1573 **Theorem 1. (Identifiability of Subspace)** Under the causal model described above, if the estimated
 1574 observations matches the true joint distribution of any $\{\mathbf{x}_{m,A}, \mathbf{x}_{m,B}, \mathbf{x}_{m,C}\}$ (they are exchangeable)
 1575 which are three measurements draw from one modality, and:

- 1576 i (Well-Posed Probability): The joint, marginal, and conditional distributions of $(\mathbf{x}_{m,B}, \mathbf{z}_m)$ are
 1577 all bounded and continuous.
- 1578 ii (Modality Variability): The operators $L_{\mathbf{x}_{m,C}|\mathbf{z}_m}$ and $L_{\mathbf{x}_{m,A}|\mathbf{x}_{m,C}}$ are injective.
- 1579 iii (Measurement Changes): For any $\mathbf{z}_t^{(1)}, \mathbf{z}_t^{(2)} \in \mathcal{Z}_t$ where $\mathbf{z}_t^{(1)} \neq \mathbf{z}_t^{(2)}$, we have $p(\mathbf{x}_{m,B}|\mathbf{z}_t^{(1)}) \neq$
 1580 $p(\mathbf{x}_{m,B}|\mathbf{z}_t^{(2)}, \mathbf{s})$.
- 1582 iv (Differentiability): There exists a functional M such that $M [p_{\mathbf{x}_{m,B}|\mathbf{z}_m, \mathbf{s}}(\cdot | \mathbf{z}_m, \mathbf{s})] = h(\mathbf{z}_m, \mathbf{s})$
 1583 for all $\mathbf{z}_m \in \mathcal{Z}_m$ and $\mathbf{s} \in \mathcal{S}$, where h is differentiable.

1584 Then we have $[\hat{\mathbf{z}}_m, \hat{\mathbf{s}}] = h(\mathbf{z}_m, \mathbf{s})$, where h is an invertible and differentiable function.
 1585

1586
 1587
 1588 **Discussion on Insufficient Measurements.** Importantly, Theorem ?? is not limited to the use of
 1589 multiple measurements within a single modality for recovering latent variables. It also reveals that,
 1590 when the number of measurements in one modality is insufficient (i.e., fewer than 3), additional
 1591 modalities can provide complementary information, provided that the required assumptions are met.

1592 We first introduce another operator to represent the point-wise distributional transformation. To
 1593 maintain generality, we denote two variables as a and b , with respective support sets \mathcal{A} and \mathcal{B} .
 1594

1595 **Definition 1. (Linear Operator) ?** Consider two random variables a and b with support \mathcal{A} and \mathcal{B} ,
 1596 the linear operator $L_{b|a}$ is defined as a mapping from a probability function p_a in some function
 1597 space $\mathcal{F}(\mathcal{A})$ onto the probability function $p_b = L_{b|a} \circ p_a$ in some function space $\mathcal{F}(\mathcal{B})$,

$$\mathcal{F}(\mathcal{A}) \rightarrow \mathcal{F}(\mathcal{B}) : p_b = L_{b|a} \circ p_a = \int_{\mathcal{A}} p_{b|a}(\cdot | a) p_a(a) da. \quad (4)$$

1602
 1603 **Definition 2. (Diagonal Operator)** Consider two random variable a and b , density functions p_a
 1604 and p_b are defined on some support \mathcal{A} and \mathcal{B} , respectively. The diagonal operator $D_{b|a}$ maps the
 1605 density function p_a to another density function $D_{b|a} \circ p_a$ defined by the pointwise multiplication of
 1606 the function $p_{b|a}$ at a fixed point b :

$$p_{b|a}(b | \cdot) p_a = D_{b|a} \circ p_a, \text{ where } D_{b|a} = p_{b|a}(b | \cdot). \quad (5)$$

1611 For brevity, we define $\mathbf{w}_m := [\mathbf{z}_m, \mathbf{s}]$ with the support set \mathcal{W}_m .
 1612
 1613
 1614
 1615
 1616

1617 *Proof.* $\mathbf{x}_{m,A}, \mathbf{x}_{m,B}, \mathbf{x}_{m,C}$ are conditional independent given \mathbf{w}_m , which implies two equations:
 1618

$$p(\mathbf{x}_{m,A} | \mathbf{x}_{m,B}, \mathbf{w}_m) = p(\mathbf{x}_{m,A} | \mathbf{w}_m), \quad p(\mathbf{x}_{m,C} | \mathbf{x}_{m,B}, \mathbf{x}_{m,A}, \mathbf{w}_m) = p(\mathbf{x}_{m,C} | \mathbf{w}_m). \quad (6)$$

We can obtain $p(\mathbf{x}_{m,C}, \mathbf{x}_{m,B} \mid \mathbf{x}_{m,A})$ directly from the observations, $p(\mathbf{x}_{m,A})$ and $p(\mathbf{x}_{m,C}, \mathbf{x}_{m,B}, \mathbf{x}_{m,A})$, and then the transformation in density function are established by

$$p(\mathbf{x}_{m,C}, \mathbf{x}_{m,B} \mid \mathbf{x}_{m,A}) = \underbrace{\int_{\mathcal{W}_m} p(\mathbf{x}_{m,C}, \mathbf{x}_{m,B}, \mathbf{w}_m \mid \mathbf{x}_{m,A}) d\mathbf{w}_m}_{\text{integration over } \mathcal{W}_m} \quad (7)$$

$$= \underbrace{\int_{\mathcal{W}_m} p(\mathbf{x}_{m,C} \mid \mathbf{x}_{m,B}, \mathbf{w}_m, \mathbf{x}_{m,A}) p(\mathbf{x}_{m,B}, \mathbf{w}_m \mid \mathbf{x}_{m,A}) d\mathbf{w}_m}_{\text{factorization of joint conditional probability}} \quad (8)$$

$$= \underbrace{\int_{\mathcal{W}_m} p(\mathbf{x}_{m,C} \mid \mathbf{w}_m) p(\mathbf{x}_{m,B}, \mathbf{w}_m \mid \mathbf{x}_{m,A}) d\mathbf{w}_m}_{\text{by } p(\mathbf{x}_{m,C} \mid \mathbf{x}_{m,B}, \mathbf{x}_{m,A}, \mathbf{w}_m) = p(\mathbf{x}_{m,C} \mid \mathbf{w}_m)} \quad (9)$$

$$= \underbrace{\int_{\mathcal{W}_m} p(\mathbf{x}_{m,C} \mid \mathbf{w}_m) p(\mathbf{x}_{m,B} \mid \mathbf{w}_m) p(\mathbf{w}_m \mid \mathbf{x}_{m,A}) d\mathbf{w}_m}_{\text{by } p(\mathbf{x}_{m,A} \mid \mathbf{x}_{m,B}, \mathbf{w}_m) = p(\mathbf{x}_{m,A} \mid \mathbf{w}_m)} \quad (10)$$

We begin by marginalizing out the variable $\mathbf{x}_{m,A}$ using the transformation structure defined in ??:

$$\begin{aligned} & \int_{\mathcal{X}_{m,A}} p(\mathbf{x}_{m,C}, \mathbf{x}_{m,B} \mid \mathbf{x}_{m,A}) p(\mathbf{x}_{m,A}) d\mathbf{x}_{m,A} = \\ & \quad \int_{\mathcal{X}_{m,A}} \int_{\mathcal{W}_m} p(\mathbf{x}_{m,C} \mid \mathbf{w}_m) p(\mathbf{x}_{m,B} \mid \mathbf{w}_m) p(\mathbf{w}_m \mid \mathbf{x}_{m,A}) p(\mathbf{x}_{m,A}) d\mathbf{w}_m d\mathbf{x}_{m,A}. \end{aligned} \quad (11)$$

This joint density expression can be rewritten using linear operators as defined in ?? and ??:

$$[L_{\mathbf{x}_{m,B}; \mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} p](\mathbf{x}_{m,C}) = [L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}} p](\mathbf{x}_{m,C}). \quad (12)$$

Thus, the composed operators satisfy the following identity:

$$L_{\mathbf{x}_{m,B}; \mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} = L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}}. \quad (13)$$

We now integrate both sides over $\mathbf{x}_{m,B} \in \mathcal{X}_{m,B}$:

$$\begin{aligned} & \int_{\mathbf{x}_{m,B} \in \mathcal{X}_{m,B}} L_{\mathbf{x}_{m,B}; \mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} d\mathbf{x}_{m,B} = \\ & \quad \int_{\mathbf{x}_{m,B} \in \mathcal{X}_{m,B}} L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}}, d\mathbf{x}_{m,B}. \end{aligned} \quad (14)$$

Since integrating out $\mathbf{x}_{m,B}$ amounts to marginalizing over the joint representation, we obtain:

$$L_{\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} = L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}}. \quad (15)$$

Assuming $L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m}$ is injective (see Assumption ??), we can invert it and obtain:

$$L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m}^{-1} L_{\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} = L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}}. \quad (16)$$

Substituting ?? into the earlier composition in ??, we derive:

$$L_{\mathbf{x}_{m,B}; \mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} = L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L_{\mathbf{w}_m \mid \mathbf{x}_{m,A}}^{-1} L_{\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}}. \quad (17)$$

Multiplying both sides of ?? by $L^{-1}\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}$ yields:

$$L_{\mathbf{x}_{m,B};\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A}} L^{-1}\mathbf{x}_{m,C} \mid \mathbf{x}_{m,A} = L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L^{-1}\mathbf{x}_{m,C} \mid \mathbf{w}_m. \quad (18)$$

The R.H.S. of ?? is in a canonical conjugation form. Under Assumption ?? and by the uniqueness of spectral decomposition (see (?), Ch. VII and (?), Theorem XV.4.5), we have:

$$\begin{aligned} L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} L^{-1}\mathbf{x}_{m,C} \mid \mathbf{w}_m &= \\ (CL_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} P)(P^{-1}D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} P)(P^{-1}L^{-1}\mathbf{x}_{m,C} \mid \mathbf{w}_m C^{-1}), \end{aligned} \quad (19)$$

where C is a nonzero scalar and P is an invertible operator encoding permutation of the eigenbasis.

This yields the following identification up to permutation and rescaling:

$$L_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} = CL_{\mathbf{x}_{m,C} \mid \hat{\mathbf{w}}_m} P, \quad D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} = P^{-1}D_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m} P. \quad (20)$$

We obtain a unique spectral decomposition in ?? with permutation and scaling indeterminacies. In the following, we will show how these indeterminacies can be resolved—if not, what informative results can still be inferred.

Since the normalizing condition

$$\int_{\mathcal{X}_{m,C}} p_{\mathbf{x}_{m,C} \mid \hat{\mathbf{w}}_m} d\mathbf{x}_{m,C} = 1 \quad (21)$$

must hold for every $\hat{\mathbf{w}}_m$, one only solution of $\int_{\mathcal{X}_{m,C}} C p_{\mathbf{x}_{m,C} \mid \mathbf{w}_m} d\mathbf{x}_{m,C} = 1$ is to set $C = 1$.

After that, we start from $D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m} = P^{-1}D_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m} P$. The operator, $D_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}$, corresponding to the set $\{p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m)\}$ for fixed $\mathbf{x}_{m,B}$ and all \mathbf{w}_m , admits a unique solution (P only change the entry position):

$$\{p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m)\} = \{p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m)\}, \quad \text{for all } \mathbf{w}_m, \hat{\mathbf{w}}_m. \quad (22)$$

Due to the set is unorder, the only way to match the R.H.S. with the L.H.S. in a consistent order is to exchange the conditioning variables, that is,

$$\{p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m^{(1)}), p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m^{(2)}), \dots\} = \quad (23)$$

$$\{p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m^{(1)}), p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m^{(2)}), \dots\} \quad (24)$$

$$\implies [p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m^{(\pi(1))}), p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m^{(\pi(2))}), \dots] = \quad (25)$$

$$[p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m^{(\pi(1))}), p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m^{(\pi(2))}), \dots] \quad (26)$$

where superscript (\cdot) denotes the index of a conditioning variable, and π is reindexing the conditioning variables. We use a relabeling map h to represent its corresponding value mapping:

$$p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid h(\mathbf{w}_m)) = p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m), \quad \text{for all } \mathbf{w}_m, \hat{\mathbf{w}}_m. \quad (27)$$

By Assumption ??, different \mathbf{w}_m corresponds to different $p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m)$, there is no repeated element in $\{p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid \mathbf{w}_m)\}$ (and $\{p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m)\}$). Hence, the relabelling map h is one-to-one (invertible).

Furthermore, Assumption 4 implies that $p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid h(\mathbf{w}_m))$ determines a unique $h(\mathbf{w}_m)$. The same holds for the $p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m)$, implying that

$$p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\mathbf{x}_{m,B} \mid h(\mathbf{w}_m)) = p_{\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m}(\mathbf{x}_{m,B} \mid \hat{\mathbf{w}}_m) \implies \hat{\mathbf{w}}_m = h(\mathbf{w}_m). \quad (28)$$

Next, Assumption ?? implies that the function h must be differentiable. Since the VAE is differentiable, we can learn a differentiable function h that satisfies Assumption ???. Consider $\hat{\mathbf{w}}_m$ related to \mathbf{w}_m via $\hat{\mathbf{w}}_m = h(\mathbf{w}_m)$. Then, we have

$$M[p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\cdot \mid \mathbf{w}_m)] = M[p_{\mathbf{x}_{m,B} \mid \mathbf{w}_m}(\cdot \mid h(\mathbf{w}_m))] = h(\mathbf{w}_m), \quad (29)$$

which is equal to $\hat{\mathbf{w}}_m$ only if h is differentiable. \square

1728 A6.3 PROOF OF THEOREM 3.2
1729

1730 Theorem ?? establishes that the modality-specific latent variables \mathbf{w}_m are block-wise identifiable.
 1731 Given multiple instances of block-wise identifiability for $[\mathbf{z}_m, \mathbf{s}]$ across different modalities m , the
 1732 shared component \mathbf{s} is expected to be identifiable as well. To support this insight, we first present a
 1733 related lemma from multi-view causal representation learning.

1734 **Lemma 1** (Identifiability from a Set of Views ?). *Consider a set of modality observations \mathbf{x}_m that
 1735 satisfy Assumption 2.1 in ?. Suppose there exists a set of modality-specific encoders, each mapping
 1736 to a common latent space. Let $\hat{g}_{\mathbf{x}_k}^{-1}$ denote a family of encoders aimed at recovering the shared
 1737 latent variables by minimizing the total entropy: $\sum_{k \in [M]} H(\hat{g}_{\mathbf{x}_k}^{-1}(\mathbf{x}_k))$. Then, under the stated
 1738 assumptions, the shared latent variables \mathbf{s} are block-identifiable.*

1740 **Theorem 2. (Identifiability of Shared Subspace)** Suppose assumptions are hold true for all the
 1741 modality and the whole latent space, and we further assume

1742 i (Entropy Regularization): $\hat{g}_{\mathbf{x}_m}^{-1}$ represent a set of shared latent variable encoders that minimizes
 1743 $\sum_{k \in [M]} H(\hat{g}_{\mathbf{x}_k}^{-1}(\mathbf{x}_k))$.

1744 Then we have the $\hat{\mathbf{s}} = h_s(\mathbf{s})$, where h_s is an invertible function.

1745
1746
1747 *Proof.* We now relate our results to Lemma ???. In ?, identifiability is established under the assumption
 1748 that multiple measurement views are available for a shared latent space, and that each measurement
 1749 process is invertible. This setting guarantees block identifiability of the latent space by aligning the
 1750 outputs of modality-specific encoders. Specifically, for each modality m , we have:

$$[\hat{\mathbf{z}}_m, \hat{\mathbf{s}}] = h(\mathbf{z}_m, \mathbf{s}), \quad (30)$$

1751 where the key insight is that $\hat{\mathbf{s}}$ corresponds to the shared component across all modality-specific
 1752 representations \mathbf{w}_m , extracted via their respective encoders.

1753 Furthermore, Lemma ?? establishes that any set of encoders minimizing the total entropy

$$\sum_{k \in [M]} H(\hat{g}_{\mathbf{x}_k}^{-1}(\mathbf{x}_k)) \quad (31)$$

1754 can recover the ground-truth shared latent variables \mathbf{s} from each modality $\mathbf{x}_m \in \mathcal{X}^m$, up to a bijective
 1755 transformation h_s :

$$\hat{\mathbf{s}} = h_s(\mathbf{s}). \quad (32)$$

1756 That is, the shared latent content \mathbf{s} is block-identified from the multi-view observations $\{\mathbf{x}_m\}_{m \in [M]}$.

1757 Finally, since each modality-specific latent variable \mathbf{z}_m is causally influenced by the shared component
 1758 \mathbf{s} , we may apply the identifiability conditions in ? as a base case. This allows us to further identify
 1759 \mathbf{z}_m up to a modality-specific bijection h_z :

$$\hat{\mathbf{z}}_m = h_z(\mathbf{z}_m). \quad (33)$$

1760 Hence, both the shared latent component \mathbf{s} and the modality-specific components \mathbf{z}_m are block-
 1761 identifiable. \square

1762
1763 **Discussion.** In the final step of our proof, we build on the identifiability result from ?, which
 1764 assumes that multiple invertible measurement processes are available to recover the shared latent
 1765 variables. In contrast, our framework relaxes this assumption by not requiring each measurement
 1766 process to be invertible. Instead, Theorem ?? ensures block identifiability of each modality-specific
 1767 latent variable \mathbf{w}_m by exploiting the information-sharing structure inherent in multi-modal and
 1768 multi-measurement settings.

1769 We further leverage a structural prior where the shared component \mathbf{s} is a common cause of the
 1770 modality-specific variables, rather than an effect. This causal asymmetry eliminates the need for
 1771 stronger conditions such as global optimization or invariance constraints. Consequently, the conditions
 1772 in ? apply, providing identifiability guarantees for the modality-specific latent variables \mathbf{z}_m .

A6.4 PROOF OF THEOREM 3.3

We begin by presenting a useful lemma from ?, which connects group-wise transformations to component-wise transformations in a Markov network. This lemma is instrumental for the subsequent proof, in particular, it enables us to first recover the latent variables within groups of adjacent nodes in the Markov network.

Lemma 2 (Identifiability of Hidden Causal Variables). *If \mathbf{z}_i is a function of at most one of $\hat{\mathbf{z}}_k$ and $\hat{\mathbf{z}}_l$, and given that \mathbf{z}_i and \mathbf{z}_j are adjacent in Markov network $\mathcal{M}_{\mathbf{z}}$, at most one of them is a function of $\hat{\mathbf{z}}_k$ or $\hat{\mathbf{z}}_l$. Then, there exists a permutation π of the estimated hidden variables, denoted as $\hat{\mathbf{z}}_{\pi}$, such that each $\hat{\mathbf{z}}_{\pi(i)}$ is a function of (a subset of) the variables in $\{\mathbf{z}_i\} \cup \Psi_{\mathbf{z}_i}$.*

Theorem 3. (Component-wise Identifiability) Suppose the assumptions (a lot abuse) in Theorem ??, Theorem ?? is satisfied, suppose we have

i (Sufficient Variability): Denote $|\mathcal{M}_{\mathbf{z}_m}|$ as the number of edges in Markov network $\mathcal{M}_{\mathbf{z}_m}$. Let

$$\begin{aligned} w(m) = & \left(\frac{\partial^3 \log p(\mathbf{z}_m | \mathbf{s})}{\partial z_{m,1}^2 \partial s_{d_s}}, \dots, \frac{\partial^3 \log p(\mathbf{z}_m | \mathbf{s})}{\partial z_{m,d_m}^2 \partial s_{d_s}} \right) \oplus \\ & \left(\frac{\partial^2 \log p(\mathbf{z}_m | \mathbf{s})}{\partial z_{m,1} \partial s_{d_s}}, \dots, \frac{\partial^2 \log p(\mathbf{z}_m | \mathbf{s})}{\partial z_{m,d_m} \partial s_{d_s}} \right) \oplus \left(\frac{\partial^3 \log p(\mathbf{z}_m | \mathbf{s})}{\partial c_{t,i} \partial c_{t,j} \partial s_{d_s}} \right)_{(i,j) \in \mathcal{E}(\mathcal{M}_{\mathbf{z}_m})}, \end{aligned} \quad (34)$$

where \oplus denotes concatenation operation and $(i, j) \in \mathcal{E}(\mathcal{M}_{\mathbf{z}_m})$ denotes all pairwise indices such that $z_{m,i}, z_{m,j}$ are adjacent in $\mathcal{M}_{\mathbf{z}_m}$. For $m \in [1, \dots, n]$, there exist $4n + |\mathcal{M}_{\mathbf{z}_m}|$ different values of \mathbf{s}_{d_s} , such that the $4n + |\mathcal{M}_{\mathbf{z}_m}|$ values of vector functions $w(m)$ are linearly independent.

ii (Sparsity Regularization): Let $\mathbf{G} \in \{0, 1\}^{d_z \times d_z}$ denote the true adjacency matrix of the latent causal graph, and $\hat{\mathbf{G}} \in \{0, 1\}^{d_z \times d_z}$ be the estimated adjacency matrix. We assume that the estimated graph is at most as dense as the true graph:

$$\|\hat{\mathbf{G}}\|_0 \leq \|\mathbf{G}\|_0,$$

where $\|\cdot\|_0$ denotes the element-wise ℓ_0 norm, i.e., the number of nonzero entries.

Then we have $\hat{\mathbf{z}}_{m,i} = h_i(\mathbf{z}_{m,\pi(j)})$, where h_i is an invertible and differentiable function.

Proof. By Theorem ??, we have

$$h(\hat{\mathbf{z}}) = \mathbf{z} \implies p_{h(\hat{\mathbf{z}})} = p_{\mathbf{z}},$$

Let J_h be the Jacobian matrix of h . The change-of-variable formula implies

$$\begin{aligned} p(\hat{\mathbf{z}} | \hat{\mathbf{s}}) |\det J_{h^{-1}}| &= p(\mathbf{z} | \mathbf{s}) \\ \log p(\hat{\mathbf{z}} | \hat{\mathbf{s}}) &= \log p(\mathbf{z} | \mathbf{s}) + \log |\det J_h|. \end{aligned} \quad (35)$$

Suppose $\hat{\mathbf{z}}_k$ and $\hat{\mathbf{z}}_l$ are conditionally independent given $\hat{\mathbf{z}}_{[n] \setminus \{k,l\}}$ i.e., they are not adjacent in the Markov network over $\hat{\mathbf{z}}$. For each $\hat{\mathbf{s}}$, by ?, we have

$$\frac{\partial^2 \log p(\hat{\mathbf{z}} | \hat{\mathbf{s}})}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l} = 0. \quad (36)$$

To see what it implies, we find the first-order derivative of Eq. equation ??:

$$\frac{\partial \log p(\hat{\mathbf{z}} | \hat{\mathbf{s}})}{\partial \hat{\mathbf{z}}_k} = \sum_{i=1}^n \frac{\partial \log p(\mathbf{z} | \mathbf{s})}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k} + \frac{\partial \log |\det J_v|}{\partial \hat{\mathbf{z}}_k}.$$

Let

$$\eta(\mathbf{s}) := \log p(\mathbf{z} | \mathbf{s}), \quad \eta'_i(\mathbf{s}) := \frac{\partial \log p(\mathbf{z} | \mathbf{s})}{\partial \mathbf{z}_i},$$

$$\eta''_{ij}(\mathbf{s}) := \frac{\partial^2 \log p(\mathbf{z} | \mathbf{s})}{\partial \mathbf{z}_i \partial \mathbf{z}_j}, \quad h'_{i,l} := \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_l}, \quad h''_{i,kl} := \frac{\partial^2 \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l}.$$

We then derive the second-order derivative w.r.t. $\hat{\mathbf{z}}_k$ and $\hat{\mathbf{z}}_l$ and apply Eq. equation ??:

$$\begin{aligned} 0 &= \sum_{j=1}^n \sum_{i=1}^n \frac{\partial^2 \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i \partial \mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \hat{\mathbf{z}}_l} \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k} + \sum_{i=1}^n \frac{\partial \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i} \frac{\partial^2 \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l} + \frac{\partial^2 \log |\det J_v|}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l} \\ &= \sum_{i=1}^n \frac{\partial^2 \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i^2} \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_l} \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k} + \sum_{j=1}^n \sum_{i: \{\mathbf{z}_j, \mathbf{z}_i\} \in \mathcal{E}(\mathcal{M}_{\mathbf{z}})} \frac{\partial^2 \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i \partial \mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \hat{\mathbf{z}}_l} \frac{\partial \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k} \\ &\quad + \sum_{i=1}^n \frac{\partial \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i} \frac{\partial^2 \mathbf{z}_i}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l} + \frac{\partial^2 \log |\det J_v|}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l} \end{aligned} \quad (37)$$

$$\begin{aligned} &= \sum_{i=1}^n \eta''_{ii}(\mathbf{s}) h'_{i,l} h'_{i,k} + \sum_{j=1}^n \sum_{i: \{\mathbf{z}_j, \mathbf{z}_i\} \in \mathcal{E}(\mathcal{M}_{\mathbf{z}})} \eta''_{ij}(\mathbf{s}) h'_{j,l} h'_{i,k} + \sum_{i=1}^n \eta'_i(\mathbf{s}) h''_{i,kl} + \frac{\partial^2 \log |\det J_v|}{\partial \hat{\mathbf{z}}_k \partial \hat{\mathbf{z}}_l}. \end{aligned} \quad (38)$$

Recall that $\mathcal{E}(\mathcal{M}_{\mathbf{z}})$ denotes the set of edges in the Markov network over Z . In the equation above, we made use of the fact that if \mathbf{z}_i and \mathbf{z}_j are not adjacent in the Markov network, then $\frac{\partial^2 \log p(\mathbf{z}|\mathbf{s})}{\partial \mathbf{z}_i \partial \mathbf{z}_j} = 0$ by ??.

By Assumption ??, consider the $2d_z + |\mathcal{M}_{\mathbf{z}}| + 1$ values of \mathbf{s} , i.e., $\mathbf{s}^{(u)}$ with $u = 0, \dots, 2d_z + |\mathcal{M}_{\mathbf{z}}|$, such that Eq. (??) hold. Then, we have $2d_z + |\mathcal{M}_{\mathbf{z}}| + 1$ such equations. Subtracting each equation corresponding to $\mathbf{s}^{(u)}$, $u = 1, \dots, 2d_z + |\mathcal{M}_{\mathbf{z}}|$ with the equation corresponding to $\mathbf{s}^{(0)}$ results in $2d_z + |\mathcal{M}_{\mathbf{z}}|$ equations:

$$\begin{aligned} 0 &= \sum_{i=1}^n (\eta''_{ii}(\mathbf{s}^{(u)}) - \eta''_{ii}(\mathbf{s}^{(0)})) h'_{i,l} h'_{i,k} + \sum_{j=1}^n \sum_{i: \{\mathbf{z}_j, \mathbf{z}_i\} \in \mathcal{E}(\mathcal{M}_{\mathbf{z}})} (\eta''_{ij}(\mathbf{s}^{(u)}) - \eta''_{ij}(\mathbf{s}^{(0)})) h'_{j,l} h'_{i,k} \\ &\quad + \sum_{i=1}^n (\eta'_i(\mathbf{s}^{(u)}) - \eta'_i(\mathbf{s}^{(0)})) h''_{i,kl}, \end{aligned}$$

where $u = 1, \dots, 2d_z + |\mathcal{M}_{\mathbf{z}}|$. Since $p_{\mathbf{z}}$ is twice continuously differentiable, we have

$$\eta''_{ij}(\mathbf{s}^{(u)}) - \eta''_{ij}(\mathbf{s}^{(0)}) = \eta''_{ji}(\mathbf{s}^{(u)}) - \eta''_{ji}(\mathbf{s}^{(0)}),$$

and therefore Eq. equation ?? can be written as

$$\begin{aligned} 0 &= \sum_{i=1}^n (\eta''_{ii}(\mathbf{s}^{(u)}) - \eta''_{ii}(\mathbf{s}^{(0)})) h'_{i,l} h'_{i,k} + \sum_{\substack{i,j: \\ i < j, \\ \{\mathbf{z}_i, \mathbf{z}_j\} \in \mathcal{E}(\mathcal{M}_{\mathbf{z}})}} (\eta''_{ij}(\mathbf{s}^{(u)}) - \eta''_{ij}(\mathbf{s}^{(0)})) (h'_{j,l} h'_{i,k} + h'_{i,l} h'_{j,k}) \\ &\quad + \sum_{i=1}^n (\eta'_i(\mathbf{s}^{(u)}) - \eta'_i(\mathbf{s}^{(0)})) h''_{i,kl}. \end{aligned}$$

Consider the vectors formed by collecting the corresponding coefficients in the equation above where $u = 1, \dots, 2d_z + |\mathcal{M}_{\mathbf{z}}|$. By Assumption A2, these $2d_z + |\mathcal{M}_{\mathbf{z}}|$ vectors are linearly independent. Thus, for any i and j such that $\{\mathbf{z}_i, \mathbf{z}_j\} \in \mathcal{E}(\mathcal{M}_{\mathbf{z}})$, we have the following equations:

$$h'_{i,k} h'_{i,l} = 0, \quad (39)$$

$$h'_{i,k} h'_{j,l} + h'_{j,k} h'_{i,l} = 0, \quad (40)$$

$$h''_{i,kl} = 0.$$

It remains to show $h'_{i,k} h'_{j,l} = 0$. Suppose by contradiction that

$$h'_{i,k} h'_{j,l} \neq 0, \quad (41)$$

which implies $h'_{i,k} \neq 0$. By Eq. equation ??, we have $h'_{i,l} = 0$, which, by plugging into Eq. equation ??, indicates $h'_{i,k} h'_{j,l} = 0$. This is a contradiction with Eq. equation ???. Thus, we must have $h'_{i,k} h'_{j,l} = 0$, which indicates that \mathbf{z}_i is a function of at most one of $\hat{\mathbf{z}}_k$ and $\hat{\mathbf{z}}_l$, and given that \mathbf{z}_i and \mathbf{z}_j are adjacent in Markov network $\mathcal{M}_{\mathbf{z}}$, at most one of them is a function of $\hat{\mathbf{z}}_k$ or $\hat{\mathbf{z}}_l$.

Then, using Lemma ??, we can obtain that there exists a permutation π of the estimated hidden variables, denoted as $\hat{\mathbf{z}}_\pi$, such that each $\hat{\mathbf{z}}_{\pi(i)}$ is a function of (a subset of) the variables in $\{\mathbf{z}_i\} \cup \Psi_{\mathbf{z}_i}$. It is worth noting that in many cases, the above result already enables us to recover some of the hidden variables up to a component-wise transformation, that is, $\hat{\mathbf{z}}_{\cdot,i} = h_i(\mathbf{z}_{\cdot,\pi(j)})$, where h_i is an invertible function. \square

We next present a proposition that shows how an arbitrary permutation over all components can be resolved into a permutation within each modality block.

Proposition 1. (*Resolving Block-Wise Permutation*) if $\hat{\mathbf{z}}_{\cdot,i} = h_i(\mathbf{z}_{\cdot,\pi(j)})$ and $\hat{\mathbf{z}}_m = h_m(\mathbf{z}_m)$ for any $m \in [M]$, we have $\hat{\mathbf{z}}_{m,i} = h_i(\mathbf{z}_{m,\pi(j)})$, where h_i is an invertible function.

Proof. Since the global mapping is given by $\hat{\mathbf{z}} = h(\mathbf{z})$, where $h = [h_1, h_2, \dots, h_M]$ acts block-wise on each modality \mathbf{z}_m , the Jacobian $J_h(\mathbf{z}) = \frac{\partial \hat{\mathbf{z}}}{\partial \mathbf{z}}$ is block-diagonal:

$$J_h(\mathbf{z}) = \begin{bmatrix} J_{h_1}(\mathbf{z}_1) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & J_{h_M}(\mathbf{z}_M) \end{bmatrix}.$$

This implies that each $\hat{\mathbf{z}}_m$ depends only on \mathbf{z}_m .

Given the global identifiability condition $\hat{\mathbf{z}}_{\cdot,i} = h_i(\mathbf{z}_{\cdot,\pi(j)})$, and the fact that both $\hat{\mathbf{z}}_{\cdot,i}$ and $\mathbf{z}_{\cdot,\pi(j)}$ must lie in the same modality m due to the block-diagonal structure, we conclude:

$$\hat{\mathbf{z}}_{m,i} = h_i(\mathbf{z}_{m,\pi(j)}).$$

\square

Discussion. We demonstrate that multi-modality information enables the use of the shared confounder \mathbf{s} as a continuous conditional prior over the modality-specific latent variables \mathbf{z}_m . This represents the key distinction from conventional multi-modality or multi-view frameworks ???. By conditioning on \mathbf{s} —for example, a gene-level representation—we can achieve component-wise identifiability of latent variables and recover their causal graph under milder assumptions. Furthermore, Proposition ?? shows that the modality-specific latent structure \mathbf{z}_m , obtained via Theorem ??, facilitates the resolution of permutation indeterminacies across the latent spaces associated with different modalities.

A7 DETAILS ABOUT NETWORK TRAINING FOR CAUSAL REPRESENTATION LEARNING

In this section, inspired by identifiability results as shown in the Theorems, we will introduce our estimation framework which enforces the proposed assumptions as constraints to identify the latent variables in each modality, in total we use several loss functions as constraints. The details are given as follows.

Network Architecture. For the high-dimensional data, we use a large foundation model to extract a high-dimensional feature first, and then use the 3-layer multi-layer perception (MLP) for the encoders and decoders. Specifically, for image data, we utilize ImageBind ? to extract 1024-dimensional embedding vectors, as this model excels at multi-modal embedding extraction. For text descriptions, we employ the gte-Qwen2-7B-instruct model from Alibaba ?, which is specifically designed for long-sentence embedding tasks and demonstrates superior performance in capturing semantic representations from extended textual content. After this gte model, we will get a 3584-dimensional embedding vector for each input text description.

Encoder and decoder. Each modality \mathbf{x}_m is given as an input to the corresponding encoder and outputs the estimated modality-specific latent $\hat{\mathbf{z}}_m$, exogenous variables $\hat{\eta}_m$, and shared latent variables \mathbf{s} across different modalities. In one modality, to ensure the conditional independence among different $\hat{\mathbf{x}}_{m,k}$ given $\hat{\mathbf{z}}_m$, $\hat{\mathbf{x}}_{m,k}$ are passed to their corresponding k -th decoders, respectively, to reconstruct

1944 the observations $\hat{\mathbf{x}}_{m,k}$ in each measurement. The reconstruction loss is calculated using the mean
 1945 squared error (MSE) as
 1946

$$\mathcal{L}_{\text{Recon}} = \sum_{m=1}^M \sum_{k=1}^{d_m} \|\mathbf{x}_{m,k} - \hat{\mathbf{x}}_{m,k}\|_2^2.$$

1951 **Conditional independence constraints.** We enforce the conditional independence condition
 1952 $\mathbf{x}_{m,j} \perp\!\!\!\perp \mathbf{x}_{m,k} \mid \mathbf{z}_m$ (where $\mathbf{x}_{m,j}$ and $\mathbf{x}_{m,k}$ refer to the j -th and k -th measurements in m -th modality)
 1953 and the independence condition on $\eta_m \perp\!\!\!\perp \mathbf{z}_m$ by enforcing the independence among components
 1954 in $\gamma = [\{\hat{\mathbf{z}}_m\}_{m=1}^M, \{\hat{\eta}_m\}_{m=1}^M, \{\hat{\epsilon}_i\}_{i=1}^{d_z}]$. To implement it, we assume that γ follows an independent
 1955 prior distribution $p(\gamma)$, such as a standard isotropic Gaussian, and enforce the independence by
 1956 matching the distribution of $\hat{\gamma}$ to the prior distribution. Specifically, we minimize the KL divergence
 1957 between the posterior and a Gaussian prior distribution as follows:

$$\mathcal{L}_{\text{Ind}} = \text{KL}(p(\gamma) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I})).$$

1960 **Proposition 2** (Conditional Independence Condition). Denote $\mathbf{x}_{m,j}$ and $\mathbf{x}_{m,k}$ are two different
 1961 measurements in one modality for the m -th modality with modality-specific latent variable \mathbf{z}_m .
 1962 $\mathbf{z}_m \subset \mathbf{z}$ is the set of block-identified latent variables, and $\eta_m \subset \eta$ are exogenous variables in
 1963 modality m . We have $\mathbf{x}_{m,j} \perp\!\!\!\perp \mathbf{x}_{m,k} \mid \mathbf{z}_m \iff \epsilon_{m,j} \perp\!\!\!\perp \epsilon_{m,k}$.

1964 **Proposition 3** (Independent Noise Condition). Denote \mathbf{z} and η as the block-identified latent variables
 1965 and exogenous variables across all modalities. ϵ 's are the causally-related noise terms. We have
 1966 $\eta \perp\!\!\!\perp \mathbf{z} \iff \eta \perp\!\!\!\perp \epsilon$.

1968 **Sparsity regularization.** We use normalization flow (?) to estimate the exogenous variables ϵ
 1969 and implement the causal relations through a learnable adjacency matrix $\hat{\mathbf{A}}$. The binary values in
 1970 $\hat{\mathbf{A}}$ represent the causal generation process between latent variables, e.g. $\hat{A}_{i,j} = 1$ indicates \hat{z}_j is the
 1971 parent of \hat{z}_i , while $\hat{A}_{i,j} = 0$ means \hat{z}_j dose not contribute to the generation of \hat{z}_i . For each component
 1972 \hat{z}_i , we select its parents $\text{Pa}(\hat{z}_i)$ based on the estimated causal adjacency matrix, and apply the flow
 1973 transformation from $\text{Pa}(\hat{z}_i)$ to $\hat{\epsilon}_i$.

1974 To encourage sparsity among the latent variables $\hat{\mathbf{z}}$, we introduce a regularization term on the learned
 1975 adjacency matrix. The sparsity assumption indicates that the optimal causal graph should be the
 1976 minimal one which still allows the model to successfully match the ground truth observational
 1977 distribution. In particular, we reduce the dependencies between different components of $\hat{\mathbf{z}}$ by adding
 1978 a \mathcal{L}_1 penalty on the adjacency matrix, s.t.,

$$\mathcal{L}_{\text{Sp}} = \|\hat{\mathbf{A}}\|_1.$$

1982 **Network Training.** In summary, the model parameters are optimized using the combination
 1983 objective:
 1984

$$\mathcal{L} = \alpha_{\text{Recon}} \mathcal{L}_{\text{Recon}} + \alpha_{\text{Ind}} \mathcal{L}_{\text{Ind}} + \alpha_{\text{Sp}} \mathcal{L}_{\text{Sp}}. \quad (42)$$

A8 DETAILS ABOUT SYNTHETIC EXPERIMENTS ON VARIANT MNIST

1989 In this section, we will introduce the synthetic experiments designed to validate our proposed
 1990 causal representation learning framework. We conduct comprehensive evaluations using carefully
 1991 constructed datasets with known causal relationships, allowing us to systematically assess the
 1992 performance of our method against established baselines.

A8.1 DETAILS ABOUT EXPERIMENTAL SETUP

1994 To systematically evaluate our proposed causal representation learning framework, we construct a
 1995 synthetic dataset with known ground-truth causal relationships using variants of the MNIST dataset.
 1996 Our synthetic dataset consists of two modalities: colored MNIST and fashion MNIST, each containing

1998 causally related latent variables. For colored MNIST, we define horizontal position as a latent cause
 1999 that influences image transparency, where digits are positioned at different horizontal locations
 2000 and their transparency varies accordingly. For fashion MNIST, we establish vertical position as
 2001 a latent cause that affects grayscale intensity of the clothing items. The causal structure connects
 2002 these modalities through a cross-modal relationship: the horizontal position in colored MNIST
 2003 serves as a causal factor for the vertical position in fashion MNIST, creating a meaningful inter-
 2004 modal dependency. Notably, our dataset design reflects different measurement characteristics across
 2005 modalities: for fashion MNIST, each sample contains a single image representing one measurement,
 2006 while for colored MNIST, we generate three images with different background colors (red, green,
 2007 blue) for each sample, providing three distinct measurements that capture different aspects of the
 2008 same underlying latent variables. The generated image examples are shown in Figure ??(a). The key
 2009 hyper-parameters are listed in Table ??.
 2010

Ground Truth Causal Graph and Training Configuration. The underlying causal relationships
 2011 in our synthetic dataset are illustrated in Figure ??(b). The causal graph demonstrates how latent
 2012 variables within and across modalities interact: horizontal position in colored MNIST causally
 2013 influences both the image transparency within the same modality and the vertical position in fashion
 2014 MNIST across modalities. Subsequently, the vertical position in fashion MNIST determines the
 2015 grayscale intensity of the fashion items. This carefully designed causal structure enables us to evaluate
 2016 whether our method can correctly identify and disentangle these known causal relationships from the
 2017 observed multi-modal data.
 2018

A8.2 DETAILS ABOUT RESULTS AND ANALYSIS

We compare our approach against several baseline methods including MCL, BetaVAE, and MMCRL
 2020 using two key metrics: R^2 (coefficient of determination) and MCC (Matthews Correlation Coefficient).
 2021 As shown in Figure Figure ??(c), our method consistently outperforms all baseline approaches across
 2022 both evaluation metrics. Specifically, our approach achieves R^2 scores of 0.96 and MCC scores
 2023 of 0.92, demonstrating superior performance in both regression and classification tasks for causal
 2024 variable identification. The substantial improvement over strong baselines like MMCRL ($R^2 = 0.90$,
 2025 $MCC = 0.85$) validates the effectiveness of our proposed framework in learning causally meaningful
 2026 representations from multi-modal observations. These results confirm that our method successfully
 2027 captures the underlying causal structure while maintaining high fidelity in representation learning,
 2028 even when dealing with asymmetric measurement structures across different modalities.
 2029

A9 DETAILS ABOUT REAL-WORLD PERSONALITY ANALYSIS ON PERSONAX

A9.1 DETAILS ABOUT EXPERIMENTAL SETUP

We conduct real-world personality analysis by training our network to extract latent representations
 2035 from both the image and text modalities of the CelebPersona dataset, followed by the application
 2036 of causal discovery to reveal underlying structures. The key hyper-parameters are listed in Table ??.
 2037 The resulting causal graph for AthlePersona is at Fig. ??(Right). For CelebPersona the causal graph
 2038 is shown in Fig. ??, we identify three shared latent variables (S_1, S_2, S_3), ten latent variables derived
 2039 from facial images ($Z_{1,1}$ to $Z_{1,10}$), and five latent variables extracted from personality descriptions
 2040 ($Z_{2,1}$ to $Z_{2,5}$). Each variable is grounded in real-world interpretable features, enabling meaningful
 2041 analysis of the causal pathways.
 2042

A9.2 DETAILS ABOUT RESULTS AND ANALYSIS

We interpret the shared latent variables S_1, S_2 , and S_3 as representing education, cultural background,
 2045 and growing environment, respectively. Notably, S_2 influences $Z_{2,4}$, which we interpret as cultural
 2046 background shaping one’s language use, while S_3 influences $Z_{2,1}$, suggesting that the growing
 2047 environment affects the use of positive language. Furthermore, expressiveness ($Z_{2,5}$) is found to
 2048 causally influence approachability ($Z_{1,10}$), reinforcing the idea that one’s ability to convey emotions
 2049 plays a key role in how approachable they appear. On the visual side, we observe that variations in
 2050 event context ($Z_{1,5}$) and lighting conditions ($Z_{1,3}$) lead to changes in hairstyle ($Z_{1,4}$), which in turn
 2051 influence face visibility ($Z_{1,8}$), overall style ($Z_{1,6}$), and how good-looking ($Z_{1,7}$) the person appears.

2052
 2053 To validate our example, we conducted an RCIT test between the Big Five traits (Final_O to Final_N)
 2054 and two sets of latent variables: five derived from personality descriptions across both datasets. We
 2055 also carry out the same tests on ten facial attributes from CelebPersona and ten latent variables derived
 2056 from facial images. As shown in Figure ?? (a), confidence ($Z_{2,1}$) exhibits strong statistical dependence
 2057 with Openness, Extraversion, and Agreeableness. In contrast, Self-awareness ($Z_{2,4}$) is significantly
 2058 associated only with Extraversion, suggesting that more extraverted individuals tend to be more
 2059 self-aware, likely due to their expressiveness, social engagement, and sensitivity.

2060 For the test result of CelebPersona ?? (b), positive language use ($Z_{2,1}$) has significant dependence
 2061 with Agreeableness indicates that more agreeable individuals are likely to use warmer and more
 2062 positive language, aligning with their prosocial and empathetic tendencies. On the other hand, the
 2063 high p-values across all Big Five traits suggest that expressiveness ($Z_{2,5}$) operates independently
 2064 of stable personality dimensions in this dataset, possibly reflecting more situational or behavioural
 2065 factors not captured by self-reported traits. In ?? (c), the p-value heatmap confirms that many facial
 2066 attributes are significantly influenced by latent appearance factors like clothing style ($Z_{1,1}$), lighting
 2067 ($Z_{1,3}$), and event context ($Z_{1,5}$), as shown in the causal graph. Traits like Big_Nose, Pointy_Nose,
 2068 and Oval_Face are tightly linked to hairstyle and good looking ($Z_{1,7}$).
 2069
 2070
 2071
 2072
 2073
 2074
 2075
 2076
 2077
 2078
 2079
 2080
 2081
 2082
 2083
 2084
 2085
 2086
 2087
 2088
 2089
 2090
 2091
 2092
 2093
 2094
 2095
 2096
 2097
 2098
 2099
 2100
 2101
 2102
 2103
 2104
 2105

Table A2: Full Table of Features and Descriptions for CelebPersona.

CelebPersona Dataset			
Feature	Type	Description	Missing Rate (%)
Id	string	Unique identifier for each celebrity	0
Height	float32	Height in centimeters	71.5
Weight	float32	Weight in kilograms	87.0
Birthday	int32	Day of birth	2.0
Birthmonth	int32	Month of birth	2.0
Birthyear	int32	Year of birth	0.6
Latitude	float32	Latitude of country's central location	0.2
Longitude	float32	Longitude of country's central location	0.2
Occupation_Num 0 = Entertainment & Performing Arts 1 = Music 2 = Sports 3 = Media & Film Production 4 = Business & Finance 5 = Academia & Science 6 = Healthcare 7 = Legal & Government 8 = Arts & Culture 9 = Religion & Service 10 = Aviation & Space 11 = Other	int32 0	[I]Occupation category:	
Gender_Num 1 = Male 2 = Female	int32 0.2	[I]Gender:	
Chatgpt_output	string	Full trait write-up by ChatGPT encoded in embeddings	0
Gemini_output	string	Full trait write-up by Gemini encoded in embeddings	0
Llama_output	string	Full trait write-up by LLaMA encoded in embeddings	0
Chatgpt_o to Chatgpt_n 0 = Unknown 1 = Disagree 2 = Neutral 3 = Agree	int32 0	[I]Big Five scores (OCEAN) by ChatGPT:	
Gemini_o to Gemini_n	int32	Big Five scores (OCEAN) by Gemini	0
Llama_o to Llama_n	int32	Big Five scores (OCEAN) by LLaMA	0
Final_o to Final_n	int32	Final aggregated scores for Big Five traits	0
Arched_Eyebrows -1 = Absent 0 = Unknown 1 = Present	int32 0	[I]Binary facial feature:	
Big_Nose	int32	Binary facial feature	0
Pointy_Nose	int32	Binary facial feature	0
Bushy_Eyebrows	int32	Binary facial feature	0
Big_Lips	int32	Binary facial feature	0
Oval_Face	int32	Binary facial feature	0
Double_Chin	int32	Binary facial feature	0
Receding_Hairline	int32	Binary facial feature	0
Narrow_Eyes	int32	Binary facial feature	0
High_Cheekbones	int32	Binary facial feature	0
Image_1 to Image_35	image	Up to 35 facial images embeddings per identity	-

2160
2161
2162
2163
2164
2165

Table A3: Summary of Terms of Use Compliance for Different Sports Leagues.

Sports Leagues Terms of Use Compliance			
Sports League	Official Website Reference	Original Statement	Requires Consent?
NBA	Terms of Use <i>§ 9. NBA STATISTICS</i> ?	"By using such NBA Statistics, you agree that: (i) any use, display, or publication of the NBA Statistics shall include a prominent attribution to NBA.com in connection with such use, display, or publication; (ii) the NBA Statistics may only be used, displayed, or published for legitimate news reporting or private, non-commercial purposes;..."	No
NFL	Terms and Conditions <i>§ 1. INTRODUCTION; GENERAL; OWNERSHIP; PROHIBITIONS</i> ?	"You may use the Services solely for your own individual non-commercial and informational purposes only. Any other use, including for any commercial purposes, is strictly prohibited without our express prior written consent."	No
MLB	Terms of Use Agreement ?	"... you must not reproduce, prepare derivative works based upon, distribute, perform or display the MLB Digital Properties without first obtaining the written permission of MLB or otherwise as expressly set forth in the terms and conditions of the MLB Digital Properties. The MLB Digital Properties must not be used in any unauthorized manner."	Yes
NHL	Terms of Service <i>§ 7. Intellectual Property</i> ?	"You may access, use, and display the Services, but only for non-commercial, informational, personal use, without modification or alteration in any way, and only so long as you comply with these Terms."	No
Premier League	Terms of Use <i>§ 6. Intellectual Property Rights</i> ?	"You may download and print material from the Website or App as is reasonable for your own private and personal use; You may also forward such material from the Website or App to other people for their private and personal use provided you credit us as its source and add the Website address."	No
La Liga	Legal Notice and Conditions of Use <i>§ 3. Use of the Website?</i>	"... The User undertakes to refrain from (a) using the Contents in a manner... (b) reproduce or copy, distribute, allow public access through any form of public communication, adapt, transform or modify the Contents, unless authorised by the owner of the corresponding rights or it is legally permitted..."	Yes
Serie A	General Terms and Conditions of the License Agreement <i>§ 2. Right limitations § 2.2 (ii) Official Data?</i>	"Except in the case of a separate written agreement between Lega Serie A and the Licensee establishing otherwise, the Licensee may only exploit the data related to the Competitions, the Matches, the Clubs and the players in the context ..."	Yes
Bundesliga	Terms of Use Services <i>§ 8. Audiovisual Content</i> ?	"The audiovisual content available within the Products is made available to the User for personal and non-commercial purposes only. The User is authorized to use this audiovisual content only for the purposes of information and entertainment in the private sphere for themselves and persons personally associated with them (e.g. family members, friends and acquaintances). Limited to these purposes, the DFL grants the User a non-exclusive, non-transferable, non-sub-licensable right of use to access and view the audiovisual content within the Products. With the exception of the aforementioned limited right of use, the User is not granted any rights to the audiovisual content."	No
Ligue 1	Terms and Conditions of Use <i>§ 6. Intellectual Property</i> ?	"... Any total or partial reproduction of the Site or its elements without prior written authorization from the publisher (LFP) may lead to legal proceedings against the infringers."	Yes
PGA Tour	Terms of Use <i>§ 7. Conduct(D)</i> ?	"You may use real time scoring, statistics and other data (whether current or archival) collected from PGATOUR.COM solely for legitimate news reporting and for personal, non-commercial purposes. You shall not use real time scoring, statistics or other data (whether current or archival) collected from PGATOUR.COM for sale, license or other commercial purposes (including, without limitation, commercial gambling purposes), unless expressly licensed by the PGA TOUR Parties."	No
ATP Tour	Terms & Conditions <i>§ 7. PROHIBITED USES</i> <i>§ 8. A. Ownership</i> ?	"ATP owns or has the right to use all of the data, information, text, images, streaming media, video, sounds, icons, scores, rankings, statistics, and other content contained on this Website (the "Content"), and the copyrights and other intellectual property rights therein, unless otherwise noted. You may print one copy of the Content of this Website for your own personal, non-commercial use."	No

2209
2210
2211
2212
2213

2214

2215

2216

Table A4: AI Model Arena Scores and API Pricing recorded on April 10 2025.

2217

2218

Model Name	Company	Arena Score	API Price (I/O)	Used by Us
Gemini-2.5-Pro-Exp-03-25	Google	1439	\$1.25/\$10.00	yes
Llama-4-Maverick-03-26-Experimental	Meta	1417	\$5.00/\$15.00	yes
ChatGPT-4o-latest (2025-03-26)	OpenAI	1410	\$2.50/\$10.00	yes
Grok-3-Preview-02-24	xAI	1403	\$3.00/\$15.00	yes
GPT-4.5-Preview	OpenAI	1398	\$75.00/\$150.00	no
Gemini-2.0-Flash-Thinking-Exp-01-21	Google	1380	\$0.10/\$0.40	yes
Gemini-2.0-Pro-Exp-02-05	Google	1380	\$0.10/\$0.40	no
DeepSeek-V3-0324	DeepSeek	1369	\$0.07/\$1.10	yes
DeepSeek-R1	DeepSeek	1358	\$0.14/\$2.19	yes
Gemini-2.0-Flash-001	Google	1354	\$0.10/\$0.40	yes
Qwen2.5-Max	Alibaba	1340	\$1.60/\$6.40	yes
QwQ-32B	Alibaba	1315	\$0.29/\$0.39	yes

2229

2230

2231

2232

2233

Table A5: Descriptions and suitability of different independence test methods used in the paper.

2234

Test	Full Name	Description	Variable Type
CSQ	Chi-Square Test	A classical test that evaluates whether two categorical variables are statistically independent.	Categorical
GSQ	G-Square Test	A likelihood-ratio version of the Chi-Square test, more robust in some small sample cases.	Categorical
RCIT	Randomized Conditional Independence Test	A non-parametric method using randomized Fourier features to approximate kernel-based CI testing.	Continuous/Mixed
HSIC	Hilbert-Schmidt Independence Criterion	A kernel-based method for measuring dependence in high-dimensional data using reproducing kernel Hilbert spaces.	Continuous/Mixed
KCI	Kernel-based Conditional Independence Test	A kernel-based extension of HSIC for testing conditional independence, suitable for complex data.	Continuous/Mixed

2256

2257

2258

2259

2260

Table A6: Key hyperparameters used in experiments.

2261

2262

2263

2264

2265

2266

2267

Hyperparameter	MNIST	PersonaX
Learning Rate	2e-6	3e-4
Training Epochs	3000	3000
Reconstruction Loss Coefficient	2	1
Conditional Independence Loss Coefficient	1e-2	1e-2
Sparsity Loss Coefficient	1e-3	1e-3

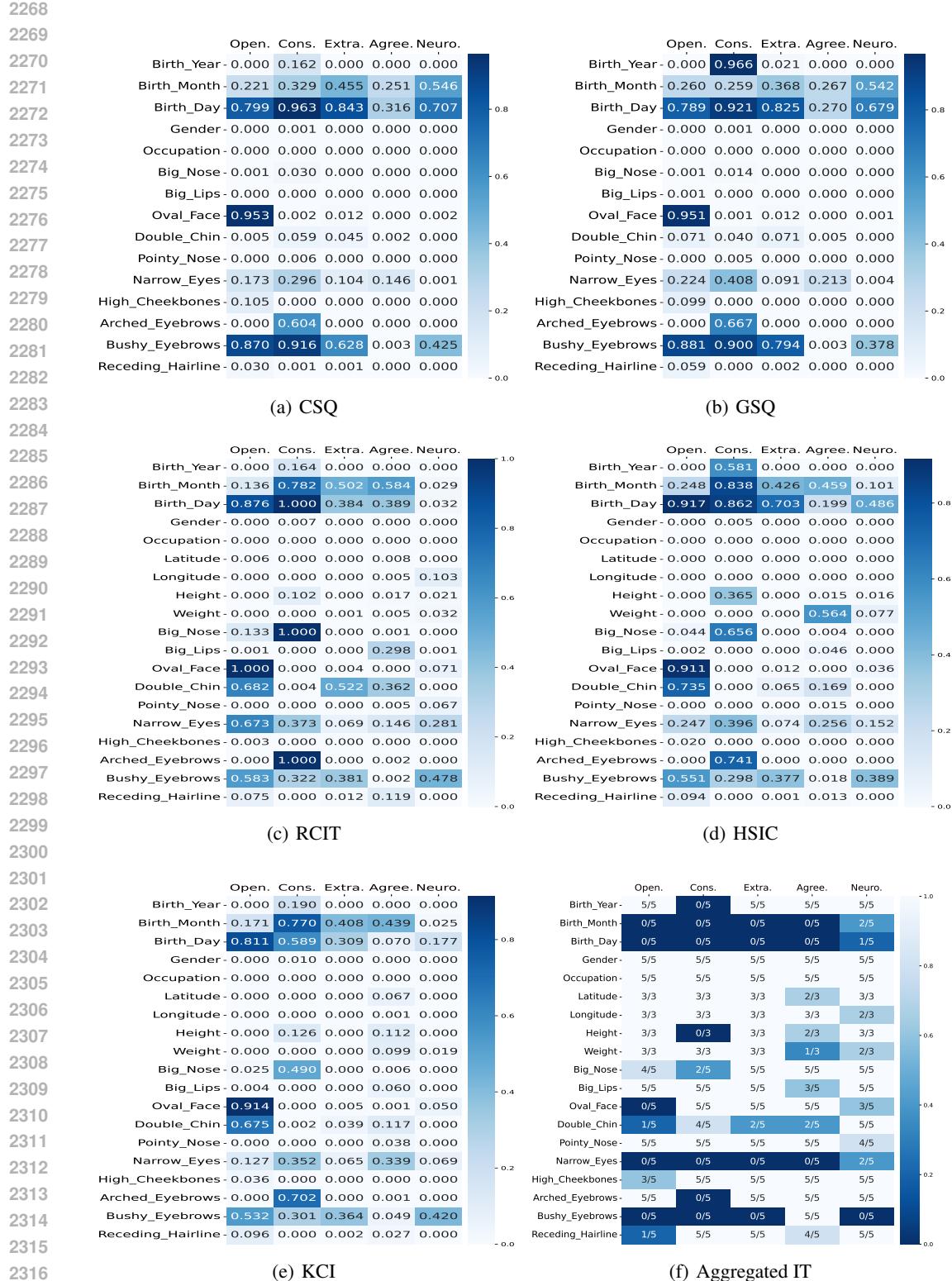


Figure A6: CelebPersona: Heatmap of P-value obtained from different independence test. (a) Chi-Square Test (b) G-Square Test (c) Randomized Conditional Independent Test (d) Hilbert-Schmidt Independence Test (e) Kernel-based Conditional Independence Test (f) Aggregated results over all above 5 methods.

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333

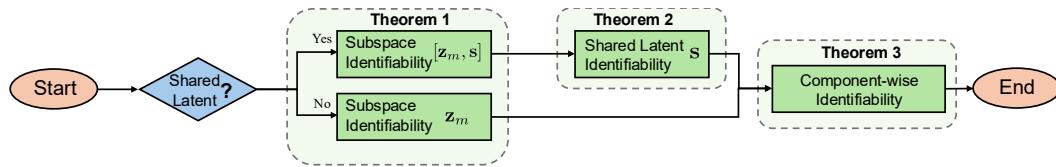


Figure A7: The high-level flowchart of the our theorems.

2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348

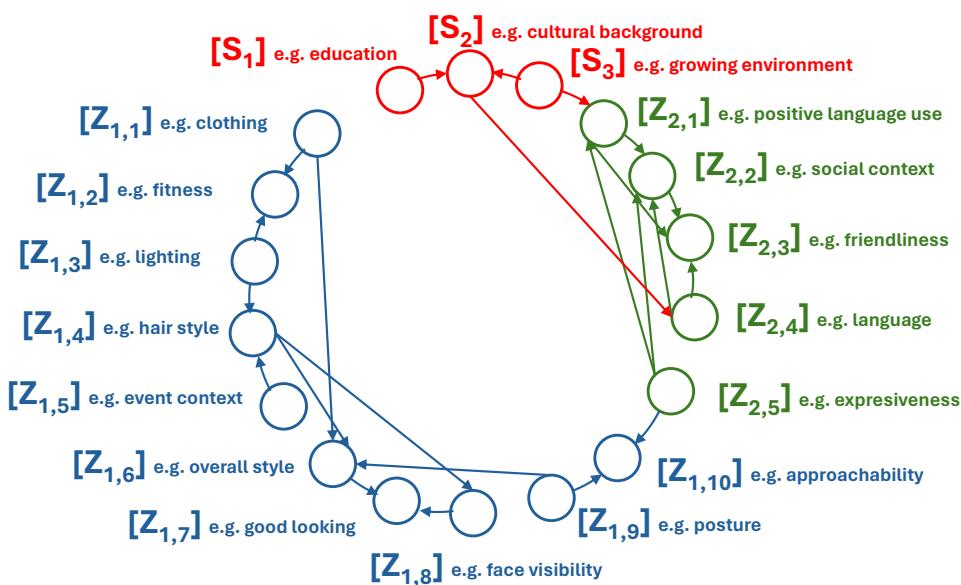


Figure A8: The causal graph with latent variables learned from CelebPersona dataset. Red, blue, and green nodes correspond to shared latents, facial image latents, and personality text latents.

2368
2369
2370
2371
2372
2373
2374
2375

