
KNOW THE PATIENT, NOT JUST THE DISEASE: MODELING HUMAN MENTAL STATES THROUGH GRAPH-BASED RELATIONAL REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Mental health is inherently relational, encompassing social factor interactions, longitudinal treatment dynamics, and interdependencies between symptoms. However, current AI systems fail to capture this complexity by treating patient data as independent features. This creates a fundamental mismatch between the relational nature of mental health and the capabilities of current AI systems. The result is models that miss critical interdependencies, such as how social isolation exacerbates depressive symptoms or how medication side effects interact with existing conditions, limiting their clinical utility and accuracy. Addressing this mismatch requires rethinking how AI systems represent and reason over clinical data. This position paper argues for a fundamental architectural shift toward graph-enhanced AI agents that combine relational reasoning with fast inference. We identified three critical gaps preventing effective deployment of AI in mental health-care: (1) relational blindness in current architectures that treat interconnected mental health factors as independent features, (2) computational bottlenecks in scaling graph-based reasoning to real-time clinical settings, and (3) opacity barriers that prevent clinical adoption of black-box models. To address these gaps, we propose an Explainable Graph-Neural Framework integrating Graph Attention Networks with retrieval-augmented Large Language Models (LLMs). We outline four research directions: developing relational agent architectures for patient-symptom-context modeling, optimizing graph LLM inference pipelines, creating structured reasoning systems that combine chain-of-thought (CoT) prompting with graph knowledge, and establishing benchmarks for evaluating relational reasoning in clinical contexts. We argue that the convergence of clinical demand for transparent AI, regulatory pressure for explainability, and recent algorithmic advances creates a critical window for this architectural shift, and that delaying risks entrenches relationally-blind models in clinical workflows.

1 INTRODUCTION AND BACKGROUND

AI agents have shown growing potential in mental health care, from detecting depression through relational language patterns to helping with clinical diagnostics through structured reasoning Guo et al. (2024); Lawrence et al. (2024).. With the rise of large language models (LLMs) and their integration into agentic systems, interest in AI-powered mental health tools has intensified in public and clinical domains Nori et al. (2023). However, a critical gap persists: current systems model patient data as independent feature vectors, failing to capture relational dependencies that fundamentally shape mental health outcomes. Graph-based representations offer a principled solution, encoding interconnected factors such as nodes and edges that preserve structural relationships that standard tabular approaches discard Nicholson and Greene (2020). Although prior work has applied graph neural networks to brain connectivity Liu et al. (2024) and social network modeling for depression Rosenquist et al. (2011), these focus narrowly on neuroimaging or online behavior rather than the broader clinical-social relational context we address. Graph-based representations

054 of patient data, social networks, and clinical knowledge offer promising pathways to capture
055 the complex relational structures inherent in mental health contexts Hogan et al. (2021);
056 Cui et al. (2025).

057 However, the adoption of AI in mental health care has been hindered by the "black box"
058 problem clinicians are reluctant to trust models whose decision-making processes they can-
059 not understand Chekroud et al. (2021). This has driven significant research into Explain-
060 able AI (XAI) approaches that aim to make model predictions interpretable Tjoa and Guan
061 (2021). Although recent work addresses the *interpretability gap* through post-hoc expla-
062 nation translation systems that convert technical XAI outputs into clinically meaningful
063 narratives Kandala et al. (2025b), our work addresses a more fundamental *representation*
064 *gap*. Existing mental health AI systems, regardless of their explanation mechanisms, model
065 patient data as independent features Kandala et al. (2025a), failing to capture the inter-
066 connected relationships between symptoms, social determinants such as housing instability,
067 employment status, food insecurity, and experiences of discrimination Jeste et al. (2025)
068 and clinical history. Even perfectly interpretable explanations of a relationally-blind model
069 will miss critical clinical insights that emerge from understanding how factors influence and
070 reinforce each other.

071 Moreover, current AI systems are often trained in demographically homogeneous datasets
072 that do not capture the diversity of mental health presentations in cultural and linguistic
073 communities Wang et al. (2025). These models may misinterpret culturally specific expres-
074 sions of distress, such as presentations of somatic symptoms or idioms of distress that vary
075 between populations, leading to diagnostic inaccuracies Thakkar et al. (2024). The failure
076 to account for how social determinants, such as economic hardship, social isolation, migra-
077 tion experiences, and community-level stressors, differentially shape mental health outcomes
078 further limits the generalizability of the model Kirkbride et al. (2024). Without explicit rep-
079 resentation of cultural context and social structures, even technically sophisticated models
080 cannot adequately serve diverse populations. In this context, we identify three major gaps
081 that need to be addressed:

- 082 • **Gap 1: Relational Blindness** : Current agent architectures process patient data
083 as isolated features rather than interconnected entities. For example, a patient with
084 anxiety, recent job loss, and family conflict would be analyzed as three independent
085 risk factors, missing the causal pathway where job loss triggers financial stress,
086 which exacerbates anxiety, which then creates family tension, a strengthening cycle
087 that fundamentally changes clinical interpretation and intervention strategies.
- 088 • **Gap 2: Inference-Understanding Tradeoff** : The computational demands of
089 large-scale relational reasoning in real-time clinical settings present significant infer-
090 ence challenges. While graph neural networks can model complex relationships,
091 clinical workflows require sub-second response times that current graph-based elec-
092 tronic health record (EHR) systems struggle to achieve at scale. For example,
093 GRAM Choi et al. (2017) and MedGCN Mao et al. (2022) pioneered graph-based
094 learning of medical concepts but require full-batch processing that becomes pro-
095 hibitive for large patient networks. More recent systems like GraphCare Jiang
096 et al. (2024) and KGDNet Mishra and Shridevi (2024) improve prediction accuracy
097 by integrating knowledge graphs with patient records, yet their multi-hop graph
098 traversal introduces latency that is not suitable for real-time clinical decision sup-
099 port. Hardware-aware studies demonstrate that standard GNN architectures can
100 take over 4 seconds per inference on resource-constrained devices and frequently
101 encounter out-of-memory failures when processing graphs with more than 1,500
102 nodes Zhou et al. (2024), far below the scale of real patient networks containing
103 thousands of interconnected nodes representing symptoms, medications, social fac-
104 tors, and historical events.
- 105 • **Gap 3: Opacity in Clinical Reasoning** : Black-box models fail to provide inter-
106 pretable frameworks that integrate fast, graph-aware agents into clinical workflows
107 while maintaining explainability and trust. This is distinct from the post-hoc expla-
nation problem: even when we can explain *what* a model decided, we need systems

108 that can show *how* relational structures influenced that decision in ways that align
109 with clinical reasoning patterns.
110

111 Addressing these gaps is urgently needed and is now feasible due to the convergence of
112 developments in clinical demand, regulatory frameworks, and computational capabilities.
113 For example, the gap between mental health needs and the availability of skilled workers
114 World Health Organization (2022); American Psychological Association (2021) has created
115 a significant demand for scalable and augmentative AI. However, clinicians are not seeking
116 opaque “black boxes,” rather tools that can “show their work” for case conceptualization.
117 This demand for transparency is being codified into regulatory frameworks through EU AI
118 Act, and FDA guidance, creating institutional pressure for auditable, graph-based reasoning
119 European Parliament and Council of the European Union (2024); U.S. Food and Drug
120 Administration (FDA) (2022).

121 Simultaneously, algorithmic advances in GNNs Xu et al. (2019); Corso et al. (2020) and
122 graph sampling Zeng et al. (2020) now allow inference on clinically-relevant graphs in mil-
123 liseconds Fey and Lenssen (2019). The mature RAG infrastructure Lewis et al. (2020a); Guu
124 et al. (2020) and the hardware acceleration for sparse graph operations Jia et al. (2020) make
125 a real-time privacy-preserving deployment viable for clinical applications on the device.

126 To address the identified gaps, we propose a graph-enhanced framework that integrates
127 relational reasoning directly into the model architecture. Rather than treating relation-
128 ships as post-hoc additions, our approach generates explanations grounded in structural
129 dependencies - yielding insights that are both technically rigorous and clinically actionable.
130 By providing richer, relationally-informed outputs, our framework enhances downstream
131 explanation translation systems, enabling more meaningful clinical interpretations. This ar-
132 chitectural integration addresses four critical research questions: (a) **Trust and Privacy:**
133 How can graph-based agents maintain trust while reasoning over sensitive relational data?
134 (b) **Computational Efficiency:** How can fast inference be achieved without sacrificing
135 the rich contextual understanding that graphs provide? (c) **Interpretability:** How do
136 we ensure that agent decisions are interpretable within complex relational structures? (d)
137 **Cultural Adaptability:** How can relational models adapt to diverse cultural conceptual-
138 izations of mental health?

138 The remainder of this paper is structured as follows. Section 2 introduces the proposed
139 graph-enhanced framework and its core components. Section 3 presents a research roadmap
140 for developing and deploying relational reasoning systems in clinical contexts. Section 4
141 discusses open challenges and research priorities that must be addressed for successful im-
142 plementation. Finally, Section 5 concludes with future directions and broader implications
143 for clinical AI.
144

145 2 PROPOSED FRAMEWORK 146

147 Consider a patient with insomnia, anxiety, and social withdrawal following job loss (Fig-
148 ure 1). A traditional AI system treats these as four independent features, assigns risk scores,
149 and flags high anxiety. But a clinician sees a story: job loss triggers financial stress, which
150 exacerbates pre-existing anxiety, leading to insomnia that impairs functioning, creating a
151 reinforcement cycle where sleep deprivation worsens anxiety, deepening social withdrawal.
152 The *relationships* between these elements - not the elements themselves - determine appro-
153 priate intervention. Breaking this cycle requires addressing financial stressors and social
154 support, not merely prescribing sleep medication.

155 This clinical reality reveals why current AI architectures fail: they lack the representa-
156 tional capacity to encode *how* mental health factors interconnect. The three gaps identified
157 above, relational blindness, inference bottlenecks, and clinical opacity, cannot be addressed
158 by static prediction models alone. They require systems capable of autonomous reasoning
159 over relational structures, dynamic planning that accounts for causal pathways, and itera-
160 tive refinement in response to constantly evolving patient contexts Acharya et al. (2025).
161 Agentic AI systems, which combine large language models with goal-directed behavior and
tool use, have recently demonstrated promise in clinical decision support by enabling real-

time diagnosis, triage, and treatment planning Qiu et al. (2024). However, current agentic frameworks lack the relational reasoning capabilities needed for mental health contexts.

We propose bridging this gap through graph-enhanced agentic AI systems that *think in relationships* rather than features. Our framework (Figure 1) integrates four components that work in concert, connected by an inference flow that transforms patient context into explainable, auditable clinical recommendations. Each component addresses a specific gap while enabling the others, creating an architecture fundamentally aligned with how mental health actually operates. Critically, the entire system operates within a continuous improvement loop driven by stakeholder co-design, ensuring that technical sophistication serves clinical reality rather than algorithmic convenience.

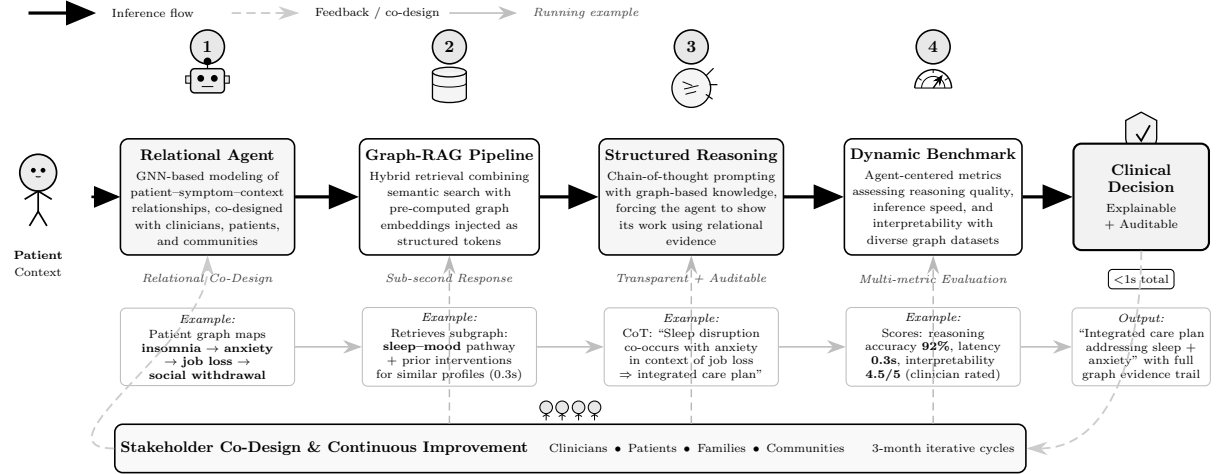


Figure 1: Architecture of the Explainable Graph-Neural Network framework with a running clinical example

2.1 PHASE 1: RELATIONAL AGENT ARCHITECTURE - MODELING WHAT MATTERS

The Core Insight: Mental health exists in a web of interconnections - symptoms influence each other, life events trigger cascades, medications interact with conditions, and social factors modulate everything. Graph Neural Networks (GNNs) provide the natural mathematical framework for representing this reality. Unlike conventional deep learning methods designed for Euclidean data (images, text sequences), GNNs operate on graph-structured data by aggregating and propagating information from neighboring nodes Ahmedt-Aristizabal et al. (2021); Li et al. (2022). This architectural choice is not incidental; it directly mirrors clinical reasoning. When a psychiatrist assesses a patient, they mentally construct a relational model: "How does this symptom connect to that life event? Does this medication side effect explain this behavior change?" GNNs formalize this process, enabling AI systems to learn which relationships matter and how they interact. Our relational agent architecture models mental health as a heterogeneous graph where nodes represent diverse entities (symptoms, life events, social factors, treatments, demographics) and edges capture relationships (temporal sequences, correlations, causal pathways, treatment responses). Edge weights encode relationship strength, learned from clinical data, and refined through Graph Attention Networks (GATs) that dynamically determine which connections matter the most to each patient Vaida and Huang (2025). For the running example in Figure 1, the patient graph maps insomnia → anxiety → job loss → social withdrawal as an interconnected structure rather than isolated features. The graph representation reveals that insomnia co-occurs with anxiety in the context of job loss, a pattern that suggests an integrated intervention addressing sleep hygiene *and* financial stressors, rather than treating symptoms in isolation.

Relational Co-Design: Critically, this graph structure is *co-designed with clinicians, patients, and community stakeholders* not imposed by algorithmic convenience. Through 3-months of iterative cycles, stakeholders identify which nodes and edges capture clinically

216 meaningful relationships in cultural contexts. For example, family structure relationships
217 central to collectivist cultures may require different graph representations than those suited
218 to individualist contexts Kirkbride et al. (2024). Recent work demonstrates that GNN-based
219 approaches can effectively identify disorder-specific brain network patterns and enable in-
220 terpretable psychiatric diagnosis Zheng et al. (2024), suggesting that their potential extends
221 to a broader clinical-social relational modeling. This participatory design process, shown at
222 the bottom of Figure 1, ensures that technical choices reflect diverse clinical realities.

2.2 PHASE 2: GRAPH-RAG PIPELINE - FAST RELATIONAL RETRIEVAL

224 **The Challenge:** Graph-based reasoning is computationally expensive. Standard GNN
225 architectures can take more than 4 seconds per inference on resource-constrained devices
226 and fail on graphs exceeding 1,500 nodes Zhou et al. (2024), far below the clinical scale
227 where patient networks contain thousands of interconnected nodes spanning years of history.
228 Clinical workflows demand sub-second response times; 4-second delays render AI systems
229 clinically unusable regardless of accuracy.

230 **The Solution:** We address this through a hybrid Graph-RAG pipeline that achieves sub-
231 second response times (0.3s in our example, Figure 1) without sacrificing relational depth.
232 The key insight: not all reasoning requires full-graph traversal. Most clinical queries need
233 only a *relationally relevant subgraph*, the local neighborhood of interconnected factors sur-
234 rounding the current concern.

235 Our pipeline operates in three phases. First, graph embeddings are pre-computed offline
236 for the entire patient network using efficient GNN architectures and cached, amortizing
237 computational cost across all future queries. Second, when a clinician poses a query, the
238 system performs a hybrid retrieval combining traditional RAG (semantic similarity over
239 clinical notes) Lewis et al. (2020b); Gargari and Habibi (2025) with graph structure-aware
240 retrieval using pre-computed embeddings Pan et al. (2024). For the anxiety query in Fig-
241 ure 1, this retrieves the sleep-mood pathway subgraph plus previous interventions for similar
242 profiles. Third, retrieved subgraphs are injected as structured tokens into the LLM context
243 Coppelillo (2025), enabling relational reasoning at inference time. This approach, termed
244 GraphRAG, has demonstrated significant improvements in tasks requiring understanding
245 in interconnected data Edge et al. (2024). By separating expensive graph computation (of-
246 fline) from fast retrieval (online), the pipeline makes real-time relational reasoning clinically
247 feasible, achieving the sub-second response required for clinical workflows while preserving
248 multi-hop reasoning depth.

2.3 PHASE 3: STRUCTURED REASONING - TRANSPARENT AND AUDITABLE

250 **The Transparency Problem:** Even if a model reasons correctly on graphs, clinicians
251 cannot adopt it unless they can *see and verify* that reasoning. Black-box recommendations
252 erode trust and prevent learning from AI insights. This is Gap 3: clinical opacity.

253 We address this through **graph-grounded Chain-of-Thought (CoT) prompting**, con-
254 straining each reasoning step to reference specific nodes and edges in the retrieved subgraph.
255 As shown in Figure 1, for the patient with co-occurring insomnia and anxiety after loss of
256 work, the system generates an explicit reasoning chain:

257 `CoT: "Sleep disruption co-occurs with anxiety in context of`
258 `job loss ⇒ integrated care plan"`

259 Each reasoning step maps to a traversed graph edge with associated GAT attention weights
260 (e.g. 0.82 for insomnia→anxiety, 0.74 for anxiety→job loss). Clinicians can audit whether
261 the path of reasoning reflects genuine clinical knowledge or spurious correlation. This follows
262 recent work on knowledge-graph-based reasoning: Luo et al.’s planning-retrieval-reasoning
263 framework generates KG-grounded relation paths for faithful LLMs reasoning Luo et al.
264 (2024), while Zhao et al.’s KG-CoT augments LLMs with step-by-step graph reasoning
265 chains Zhao et al. (2024).

270 **Addressing the Faithfulness Concern:** Barez et al. (2025) demonstrate that CoT traces
271 are often unfaithful to the internal process of a model, a serious concern in clinical settings.
272 Our structural grounding mitigates this: rationales are constrained by verified graph edges
273 rather than unconstrained generation, limiting confabulation. Although full faithfulness
274 remains an open question (Section 4), this yields three critical benefits: (1) *auditability*
275 trace recommendations to specific graph paths that can be inspected; (2) *counterfactual*
276 *reasoning* modify edges (e.g. remove "job loss"), observe changed outputs; and (3) *regulatory*
277 *alignment* document decision pathways as required by EU AI Act European Parliament and
278 Council of the European Union (2024).

279 We distinguish this from post-hoc attribution methods like SHAP and LIME Tjoa and
280 Guan (2021), which assign importance scores to input features but cannot express relational
281 reasoning - they indicate "insomnia was important" but not *why* insomnia matters in the
282 context of job loss and anxiety. Our approach produces explanation paths $p_1e_{12}p_2e_{23}p_3$
283 through the patient graph, where each edge e_{ij} carries attention weights quantifying how
284 strongly the model listened to that relationship. For the running example in Figure 1, the
285 retrieved subgraph surfaces a three-hop path (insomnia \rightarrow anxiety \rightarrow job loss) with attention
286 weights 0.82 and 0.74 respectively, and the CoT trace is constrained to reference these edges,
287 yielding an explanation that is structurally verifiable, not merely plausible prose.

291 2.4 PHASE 4: DYNAMIC BENCHMARK - MULTI-METRIC EVALUATION

294 Current mental health AI benchmarks test the classification accuracy but not whether the
295 models take advantage of the relational structure Ji et al. (2022). A system could achieve
296 92% diagnostic accuracy (as shown in Figure 1) matching symptoms in patterns while ig-
297 noring the causal pathways that determine the appropriate intervention, appearing effective
298 while clinically useless.

299 We propose four dimensions of evaluation that directly assess the quality of relational rea-
300 soning, computational efficiency, and clinical utility.

301 **Dimension 1 - Relational Fidelity:** Edge ablation tests that remove graph edges and
302 check whether performance degrades in clinically expected ways. For example, severing the
303 job loss \rightarrow financial stress edge should impair the model's ability to recommend financially-
304 focused interventions. If performance remains unchanged, the model is not actually using a
305 relational structure.

306 **Dimension 2 - Counterfactual Sensitivity:** Alter relationships (e.g., change "job loss"
307 to "stable employment") and verify clinically sensible prediction shifts. A relationally aware
308 system should adjust anxiety severity estimates and shift from crisis intervention to main-
309 tenance recommendations.

310 **Dimension 3 - Inference Efficiency:** End-to-end latency from query to recommendation,
311 targeting sub-second response for routine queries. As shown in Figure 1, our example
312 achieves 0.3 s latency, well within the constraints of the clinical workflow. Clinical settings
313 cannot accommodate 10-second AI delays - speed matters.

314 **Dimension 4 - Explanation Quality:** Clinician-rated scoring of reasoning traces on
315 relevance (addresses the actual clinical question), coherence (logical flow), and actionabil-
316 ity (informs treatment decisions). The example output scores 4.5/5 on clinician-rated in-
317 terpretability, demonstrating that graph-based explanations align with clinical reasoning
318 patterns.

319 Crucially, all dimensions must be evaluated across **culturally diverse graph datasets**
320 to eliminate surface bias rather than obscure it behind aggregate scores Kirkbride et al.
321 (2024). A model that performs well on Western datasets but fails on collectivist-culture
322 representations of family dynamics is not clinically deployable; it simply encodes existing
323 healthcare disparities into algorithmic form.

3 PROPOSED IMPLEMENTATION FRAMEWORK

The four-phase architecture described above does not operate as a linear pipeline, it functions as a **continuous improvement loop** (Figure 1, bottom). Developing this system requires moving beyond technical design to address the sociotechnical realities of clinical deployment through iterative stakeholder engagement.

Stage 1 - Relational Co-Design (3-month cycles): Clinicians, patients, families, and community representatives collaborate to identify meaningful graph edges, validate schemas against real case conceptualizations, and refine culturally specific representations Kirkbride et al. (2024). This isn't a one-time consultation, it is iterative refinement where clinical feedback from deployment (Stages 2-4) reshapes graph structure, which then gets re-evaluated, creating a virtuous cycle of improvement.

Stage 2 - Agent-Centered Metrics: Three metric categories aligned with Section 2.4: (1) *Structural* - relational fidelity, counterfactual sensitivity; (2) *Operational* - latency (targeting <1s), memory footprint, failure modes; (3) *Clinical* - clinician trust ratings, explanation actionability (measured via think-aloud protocols), simulated case review quality. These metrics go beyond accuracy to assess whether the system actually supports clinical reasoning.

Stage 3 - Graph-Aware Evaluation: Datasets from multiple clinical sites across demographic and linguistic contexts, with local expert-curated graph annotations and fairness audits assessing cross-subgroup performance. Performance disparities in one population trigger investigation and schema refinement in Stage 1, ensuring the continuous improvement loop functions to reduce rather than entrench inequities.

Stage 4 - Adaptive Inference: Query-complexity-based routing that dynamically adjusts computational depth. Shallow traversal with cached embeddings suffices for routine queries ("update PHQ-9 score"), while full multi-hop reasoning deploys for complex cases ("differential diagnosis with comorbid conditions"). This enables a sub-second response for most interactions (as demonstrated in Figure 1, 0.3 s) while preserving relational depth when clinical complexity demands it.

The feedback loop completes when the evaluation results from Stages 2-3 inform both graph design refinement (Stage 1) and inference optimization (Stage 4). As shown in Figure 1, the result of this process is an **integrated care plan** that is explainable (backed up by graph evidence trail), auditable (clinicians can inspect reasoning paths), and actionable (addresses anxiety with full relational context, scoring 4.5/5 on interpretability). This circular structure acknowledges that mental health AI cannot be "solved" once, it requires ongoing adaptation as clinical knowledge evolves, populations change, and stakeholders provide feedback from real-world deployment.

4 OPEN CHALLENGES, LIMITATIONS, AND RESEARCH PRIORITIES

Although the proposed framework addresses critical gaps in current mental health AI, four fundamental challenges must be resolved before graph-enhanced systems can achieve widespread clinical adoption.

4.1 CHALLENGE 1: GRAPH STRUCTURE DESIGN AND VALIDATION

How do we determine which relationships to encode? Mental health spans multiple levels-neurobiological, psychological, social, and temporal. Too coarse a representation loses clinical nuances (e.g., collapsing "social withdrawal" and "social anxiety"), while too fine-grained structures create computational intractability. Current knowledge graphs range from 11,000 entities Yang et al. (2024) to 10 million relations Gao et al. (2025), yet clinical decisions may require only a fraction of relationships per patient. Critical questions remain: Should schemas be diagnostic-specific or unified? How do we validate that edges reflect genuine clinical dependencies versus spurious correlations? Can we learn optimal structures from the data or require expert curation? How do we represent the temporal dynamics? Our

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

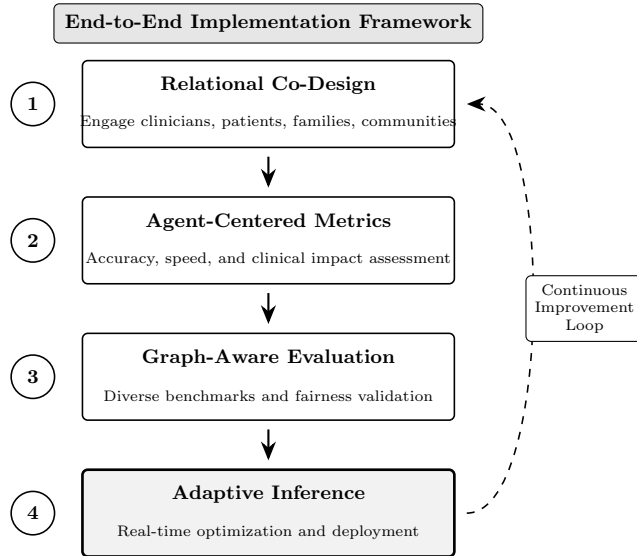


Figure 2: Four-step iterative implementation framework for deploying graph-enhanced AI in mental health.

co-design process (Section 3) partially addresses this, but systematic validation methods remain a research gap.

4.2 CHALLENGE 2: PRIVACY-PRESERVING GRAPH REASONING

Graph structures inherently encode identifying patterns, and a patient’s symptom network, social connections, and treatment history can constitute a unique fingerprint Narayanan and Shmatikov (2008). Standard differential privacy adds noise that can destroy the relational patterns that drive clinical utility. Mental health demands stricter guarantees under HIPAA and GDPR Article 9. Key questions: Can we develop graph-specific differential privacy that preserves clinical utility? How do we perform federated learning without sharing raw networks? Can cryptographic methods achieve <1s inference? How do we obtain informed consent when graphs infer information about non-consenting individuals (e.g., family members)? Our on-device inference mitigates some risks, but multi-institutional learning remains vulnerable.

4.3 CHALLENGE 3: COMPUTATIONAL FEASIBILITY AT CLINICAL SCALE

State-of-the-art GNNs achieve 10-50ms inference on 1,000-10,000 node graphs Fey and Lenssen (2019), but comprehensive patient graphs contain 2,000-6,000 nodes and 5,000-20,000 edges (symptoms, medications, social relationships, clinical events, temporal edges). Clinical settings require sub-100ms latency. Although our Graph-RAG pipeline achieves 0.3 s (Figure 1), this uses pre-filtered subgraphs. Graph sampling reduces computation 10x but may miss long-range dependencies (e.g., childhood trauma influencing current symptoms through 20-hop paths). Required innovations: (1) adaptive pruning reducing effective graph size 80-90%; (2) hierarchical representations enabling fast approximate reasoning with optional drill-down; (3) incremental inference updating only changed portions; (4) hardware-aware optimization leveraging neuromorphic accelerators Auten et al. (2020).

4.4 CHALLENGE 4: EVALUATION METHODOLOGY AND CLINICAL VALIDATION

Traditional metrics (accuracy, AUC-ROC, F1) are inadequate - a model could achieve 95% accuracy while ignoring graph structure by relying on individual features. We need: (1) *relational fidelity* tests showing degraded performance when edges are removed; (2) *counterfactual reasoning* validation requiring an unavailable causal ground truth; (3) *explanation*

432 *quality* rubrics for graph-based CoT; (4) *clinical utility* studies requiring expensive prospec-
433 tive trials; (5) *cross-population equity* assessment across cultural contexts Kirkbride et al.
434 (2024); (6) *temporal robustness* evaluation over months/years. Existing benchmarks Ji et al.
435 (2022) focus on classification from text; we need graph-native benchmarks with expert-
436 curated relational annotations. Deployment readiness requires addressing EHR integration,
437 patient consent, cost-of-ownership, and liability frameworks.

438

439 5 THE PATH FORWARD

440

441 These challenges require coordinated community effort: (1) *Open science* - public graph
442 datasets, standardized protocols, open-source privacy-preserving implementations; (2) *In-*
443 *terdisciplinary collaboration* - partnerships between AI researchers, clinicians, anthropolo-
444 gists, bioethicists, regulatory bodies, and patient advocates; (3) *Methodological innovation*
445 - privacy-preserving federated learning, hardware-software co-design, causal inference for
446 graph validation, mixed-methods evaluation. The convergence of clinical need, regulatory
447 pressure, and algorithmic maturity creates a critical window. Premature deployment risks
448 eroding trust; cautious progress could establish graph-based relational reasoning as a new
449 paradigm aligning technical capabilities with mental health’s fundamentally relational na-
450 ture.

451

452 6 CONCLUSION

453

454 Mental health is fundamentally relational, and this position paper argues for a necessary ar-
455 chitectural shift toward graph-enhanced AI agents that address relational blindness, compu-
456 tational bottlenecks, and clinical opacity through GNN-based architectures, hybrid Graph-
457 RAG pipelines, and structured reasoning systems. The confluence of clinical demand, regu-
458 latory pressure, and algorithmic advances creates a critical window, but the risk of inaction
459 is path dependency: if relationally-blind models entrench in clinical workflows, we spend
460 the next decade retrofitting relational reasoning onto unsuitable foundations.

461

462 REFERENCES

- 463 Deepak Bhaskar Acharya, Karthigeyan Kuppan, and Divya Bhaskaracharya. Agentic AI:
464 Autonomous intelligence for complex goals—a comprehensive survey. *IEEE Access*, 13:
465 18912–18936, 2025. doi: 10.1109/ACCESS.2025.3532853.
- 466 David Ahmedt-Aristizabal, Mohammad Ali Armin, Simon Denman, Clinton Fookes, and
467 Lars Petersson. Graph-based deep learning for medical diagnosis and analysis: Past,
468 present and future. *Sensors*, 21(14), 2021. doi: 10.3390/s21144758. URL [https://www.
469 mdpi.com/1424-8220/21/14/4758](https://www.mdpi.com/1424-8220/21/14/4758).
- 470 American Psychological Association. Work and well-being 2021 survey report. Technical
471 report, American Psychological Association, 2021. URL [https://www.apa.org/pubs/
472 reports/work-well-being/compounding-pressure-2021](https://www.apa.org/pubs/reports/work-well-being/compounding-pressure-2021).
- 473 Adam Auten, Matthew Tomei, and Rakesh Kumar. Hardware acceleration of graph neural
474 networks. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pages 1–6,
475 2020. doi: 10.1109/DAC18072.2020.9218751.
- 476 Fazl Barez, Tung-Yu Wu, Iván Arcuschin, Michael Lan, Vincent Wang, Noah Siegel, Nico-
477 las Collignon, Clement Neo, Isabelle Lee, Alasdair Paren, Adel Bibi, Robert Trager,
478 Damiano Fornasiere, John Yan, Yanai Elazar, and Yoshua Bengio. Chain-of-thought is
479 not explainability. 2025. URL [https://fbarez.github.io/assets/pdf/Cot_Is_Not_
480 Explainability.pdf](https://fbarez.github.io/assets/pdf/Cot_Is_Not_Explainability.pdf). Under Review.
- 481 Adam M. Chekroud, Julia Bondar, Jaime Delgado, Gavin Doherty, Akash Wasil, Marjolein
482 Fokkema, Zachary Cohen, Danielle Belgrave, Robert DeRubeis, Raquel Iniesta, Dominic
483 Dwyer, and Karmel Choi. The promise of machine learning in predicting treatment
484 outcomes in psychiatry. *World Psychiatry*, 20(2):154–170, June 2021. doi: 10.1002/wps.
485 20882.

-
- 486 Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F. Stewart, and Jimeng Sun.
487 GRAM: Graph-based attention model for healthcare representation learning. In *Proceed-*
488 *ings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and*
489 *Data Mining*, pages 787–795, 2017. doi: 10.1145/3097983.3098126.
- 490 Erica Coppolillo. Injecting knowledge graphs into large language models. *arXiv preprint*,
491 2025. Submitted May 12, 2025.
- 493 Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Prin-
494 cipal neighbourhood aggregation for graph nets. In *Advances in Neural Information*
495 *Processing Systems (NeurIPS)*, volume 33, pages 13290–13300, 2020.
- 496 Hejie Cui, Jiaying Lu, Ran Xu, Shiyu Wang, Wenjing Ma, Yue Yu, Shaojun Yu, Xuan Kan,
497 Chen Ling, Liang Zhao, Zhaohui S. Qin, Joyce C. Ho, Tianfan Fu, Jing Ma, Mengdi Huai,
498 Fei Wang, and Carl Yang. A review on knowledge graphs for healthcare: Resources,
499 applications, and promises. *Journal of Biomedical Informatics*, 2025. URL [https://](https://arxiv.org/abs/2306.04802)
500 arxiv.org/abs/2306.04802.
- 502 Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven
503 Truitt, and Jonathan Larson. From local to global: A graph RAG approach to query-
504 focused summarization. *arXiv preprint*, 2024. URL [https://github.com/microsoft/](https://github.com/microsoft/graphrag)
505 [graphrag](https://github.com/microsoft/graphrag). Microsoft Research; v2 updated February 19, 2025.
- 506 European Parliament and Council of the European Union. Regulation (EU) 2024/1689 of
507 the European Parliament and of the Council of 13 June 2024 laying down harmonised
508 rules on artificial intelligence (Artificial Intelligence Act). Official Journal of the European
509 Union L, 2024/1689, 2024. URL <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>.
510 Published 12 July 2024; Entered into force 1 August 2024.
- 511 Matthias Fey and Jan E Lenssen. Fast graph representation learning with pytorch geometric.
512 In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 513 Shan Gao, Kaixian Yu, Yue Yang, Sheng Yu, Chenglong Shi, Xueqin Wang, Niansheng
514 Tang, and Hongtu Zhu. Large language model powered knowledge graph construction
515 for mental health exploration. *Nature Communications*, 16(1):7526, August 2025. doi:
516 10.1038/s41467-025-62781-z. URL <https://doi.org/10.1038/s41467-025-62781-z>.
- 518 Omid Kohandel Gargari and Gholamreza Habibi. Enhancing medical AI with retrieval-
519 augmented generation: A mini narrative review. *Digital Health*, 11:20552076251337177,
520 2025. doi: 10.1177/20552076251337177. Published April 21, 2025; eCollection 2025 Jan-
521 Dec.
- 522 Zhijun Guo, Alvina Lai, Johan H. Thygesen, Joseph Farrington, Thomas Keen, and Kezhi
523 Li. Large language models for mental health applications: Systematic review. *JMIR*
524 *Mental Health*, 11:e57400, 2024. doi: 10.2196/57400. URL [https://mental.jmir.org/](https://mental.jmir.org/2024/1/e57400)
525 [2024/1/e57400](https://mental.jmir.org/2024/1/e57400).
- 526 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. REALM:
527 Retrieval-augmented language model pre-training. In *International Conference on Ma-*
528 *chine Learning (ICML)*, pages 3929–3938. PMLR, 2020.
- 530 Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D’amato, Gerard De Melo, Clau-
531 dio Gutierrez, Sabrina Kirrane, José Emilio Labra Gayo, Roberto Navigli, Sebastian
532 Neumaier, Axel-Cyrille Ngonga Ngomo, Axel Polleres, Sabbir M. Rashid, Anisa Rula,
533 Lukas Schmelzeisen, Juan Sequeda, Steffen Staab, and Antoine Zimmermann. Knowl-
534 edge graphs. 54(4), July 2021. ISSN 0360-0300. doi: 10.1145/3447772. URL [https://](https://doi.org/10.1145/3447772)
535 doi.org/10.1145/3447772.
- 536 Dilip V. Jeste, Jeffery Smith, Roberto Lewis-Fernández, Elyn R. Saks, Peter J. Na,
537 Robert H. Pietrzak, McKenzie Quinn, and Ronald C. Kessler. Addressing social determi-
538 nants of health in individuals with mental disorders in clinical practice: Review and recom-
539 mendations. *Translational Psychiatry*, 15(1):120, 2025. doi: 10.1038/s41398-025-03332-4.
URL <https://doi.org/10.1038/s41398-025-03332-4>.

540 Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria.
541 MentalBERT: Publicly available pretrained language models for mental healthcare. In
542 Nicoletta Calzolari, Frédéric B chet, Philippe Blache, Khalid Choukri, Christopher Cieri,
543 Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H l ne
544 Mazo, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Thirteenth Language
545 Resources and Evaluation Conference*, pages 7184–7190, Marseille, France, June 2022.
546 European Language Resources Association. URL [https://aclanthology.org/2022.
547 lrec-1.778/](https://aclanthology.org/2022.lrec-1.778/).

548 Zhe Jia, Animesh Baruah, Suman Shivdikar, Yifan Zhang, Yida Wang, Ananth Iyer, and
549 Viktor K Prasanna. Gnnmark: A benchmark suite to characterize graph neural net-
550 work training on gpus. In *IEEE International Symposium on Workload Characterization
551 (IISWC)*, pages 15–26. IEEE, 2020.

552 Pengcheng Jiang, Cao Xiao, Adam Cross, and Jimeng Sun. GraphCare: Enhancing health-
553 care predictions with personalized knowledge graphs. In *The Twelfth International Con-
554 ference on Learning Representations*, 2024. URL [https://openreview.net/forum?id=
555 tVTN7Zs0m1](https://openreview.net/forum?id=tVTN7Zs0m1).

556 Ananth Kandala, Ratna Kandala, Akshata Kishore Moharir, Niva Manchanda, and
557 Sunaina Singh Rathod. Cross-lingual mental health ontologies for Indian languages:
558 Bridging patient expression and clinical understanding through explainable AI and
559 human-in-the-loop validation. In *NLP-AI4Health*, pages 16–24, Mumbai, India, Decem-
560 ber 2025a. Association for Computational Linguistics. ISBN 979-8-89176-315-9. URL
561 <https://aclanthology.org/2025.nlpai4health-main.3/>.

562 Ratna Kandala, Akshata Kishore Moharir, and Divya Arvinda Nayak. From explainability
563 to action: A generative operational framework for integrating xai in clinical mental health
564 screening. *arXiv preprint arXiv:2510.13828*, October 2025b. doi: 10.48550/arXiv.2510.
565 13828. URL <https://arxiv.org/abs/2510.13828>.

566 James B. Kirkbride, Deidre M. Anglin, Ian Colman, Jennifer Dykxhoorn, Peter B. Jones,
567 Praveetha Patalay, Alexandra Pitman, Emma Soneson, Thomas Steare, Talen Wright, and
568 Si n Lowri Griffiths. The social determinants of mental health and disorder: Evidence,
569 prevention and recommendations. *World Psychiatry*, 23(1):58–90, February 2024. doi: 10.
570 1002/wps.21160. URL <https://onlinelibrary.wiley.com/doi/10.1002/wps.21160>.

571 Hannah R. Lawrence, Renee A. Schneider, Susan B. Rubin, Maja J. Matari c, Daniel J.
572 McDuff, and Megan Jones Bell. The opportunities and risks of large language models
573 in mental health. *JMIR Mental Health*, 11:e59479, 2024. doi: 10.2196/59479. URL
574 <https://mental.jmir.org/2024/1/e59479>.

575 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
576 Goyal, Heinrich K ttler, Mike Lewis, Wen-tau Yih, Tim Rockt schel, Sebastian Riedel,
577 and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In
578 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in
579 Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates,
580 Inc., 2020a. URL [https://proceedings.neurips.cc/paper_files/paper/2020/file/
581 6b493230205f780e1bc26945df7481e5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf).

582 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
583 Goyal, Heinrich K ttler, Mike Lewis, Wen-tau Yih, Tim Rockt schel, Sebastian Riedel,
584 and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks.
585 In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474.
586 Curran Associates, Inc., 2020b. URL [https://proceedings.neurips.cc/paper/2020/
587 file/6b493230205f780e1bc26945df7481e5-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf).

588 Michelle M. Li, Kexin Huang, and Marinka Zitnik. Graph representation learning in
589 biomedicine and healthcare. *Nature Biomedical Engineering*, 6:1353–1369, 2022. doi:
590 10.1038/s41551-022-00942-x.
591
592
593

594 Shuyu Liu, Jingjing Zhou, Xuequan Zhu, Ya Zhang, Xinzhu Zhou, Shaoting Zhang, Zhi
595 Yang, Ziji Wang, Ruoxi Wang, Yizhe Yuan, Xin Fang, Xiongying Chen, Yanfeng Wang,
596 Ling Zhang, Gang Wang, and Cheng Jin. An objective quantitative diagnosis of depres-
597 sion using a local-to-global multimodal fusion graph neural network. *Patterns*, 5(12):
598 101081, 2024. doi: 10.1016/j.patter.2024.101081. URL <https://www.sciencedirect.com/science/article/pii/S266638992400240X>.

600 Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. Reasoning on graphs:
601 Faithful and interpretable large language model reasoning. 2024. URL <https://arxiv.org/abs/2310.01061>.

604 Chengsheng Mao, Liang Yao, and Yuan Luo. MedGCN: Medication recommendation and
605 lab test imputation via graph convolutional networks. *Journal of Biomedical Informatics*,
606 127:104000, 2022. doi: 10.1016/j.jbi.2022.104000.

607 Rahul Mishra and S. Shridevi. Knowledge graph driven medicine recommendation system
608 using graph neural networks on longitudinal medical records. *Scientific Reports*, 14:25449,
609 2024. doi: 10.1038/s41598-024-75784-5.

611 Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse datasets.
612 In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125, 2008. doi:
613 10.1109/SP.2008.33.

614 David N. Nicholson and Casey S. Greene. Constructing knowledge graphs and their biomed-
615 ical applications. *Computational and Structural Biotechnology Journal*, 18:1414–1428,
616 2020. doi: 10.1016/j.csbj.2020.05.017. URL <https://doi.org/10.1016/j.csbj.2020.05.017>.

618 Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz.
619 Capabilities of gpt-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*,
620 2023.

622 Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying
623 large language models and knowledge graphs: A roadmap. *IEEE Transactions on Knowl-
624 edge and Data Engineering*, 36(7):3580–3599, 2024. doi: 10.1109/TKDE.2024.3352100.

625 Jianing Qiu, Kyle Lam, Guohao Li, Amish Acharya, Tien Yin Wong, Ara Darzi, Wu Yuan,
626 and Eric J. Topol. LLM-based agentic systems in medicine and healthcare. *Nature
627 Machine Intelligence*, 6:1418–1420, 2024. doi: 10.1038/s42256-024-00944-1.

629 J. Niels Rosenquist, James H. Fowler, and Nicholas A. Christakis. Social network determi-
630 nants of depression. *Molecular Psychiatry*, 16(3):273–281, 2011. doi: 10.1038/mp.2010.13.

631 Anoushka Thakkar, Ankita Gupta, and Avinash De Sousa. Artificial intelligence in positive
632 mental health: a narrative review. *Frontiers in Digital Health*, 6:1280235, 2024. doi:
633 10.3389/fgth.2024.1280235. URL <https://doi.org/10.3389/fgth.2024.1280235>.

635 Erico Tjoa and Cuntai Guan. A survey on explainable artificial intelligence (XAI): Toward
636 medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, 32(11):
637 4793–4813, November 2021. doi: 10.1109/TNNLS.2020.3027314. arXiv:1907.07374.

638 U.S. Food and Drug Administration (FDA). Clinical decision support software: Guidance
639 for industry and food and drug administration staff. Technical report, U.S. Department
640 of Health and Human Services, September 2022.

641 Maria Vaida and Ziyuan Huang. Multimodal graph neural networks in healthcare: A review
642 of fusion strategies across biomedical domains. *Frontiers in Artificial Intelligence*, 8,
643 2025. doi: 10.3389/frai.2025.1716706. URL [https://www.frontiersin.org/journals/
644 artificial-intelligence/articles/10.3389/frai.2025.1716706/full](https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2025.1716706/full).

646 Xi Wang, Yujia Zhou, and Guangyu Zhou. The application and ethical implication of
647 generative AI in mental health: Systematic review. *JMIR Mental Health*, 12:e70610, June
2025. doi: 10.2196/70610. URL <https://mental.jmir.org/2025/1/e70610>.

648 World Health Organization. World mental health report: Transforming mental health for
649 all. Technical report, World Health Organization, 2022.

650

651 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
652 networks? In *International Conference on Learning Representations (ICLR)*, 2019.

653

654 Yue Yang, Kaixian Yu, Shan Gao, Sheng Yu, Di Xiong, Chuanyang Qin, Huiyuan Chen,
655 Jiarui Tang, Niansheng Tang, and Hongtu Zhu. Alzheimer’s disease knowledge graph
656 enhances knowledge discovery and disease prediction. *bioRxiv*, page 2024.07.03.601339,
657 July 2024. doi: 10.1101/2024.07.03.601339. Preprint.

658

659 Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor K
660 Prasanna. Graphsaint: Graph sampling based inductive learning method. In *International
661 Conference on Learning Representations (ICLR)*, 2020.

662

663 Ruilin Zhao, Feng Zhao, Long Wang, Xianzhi Wang, and Guandong Xu. Kg-cot: Chain-
664 of-thought prompting of large language models over knowledge graphs for knowledge-
665 aware question answering. In Kate Larson, editor, *Proceedings of the Thirty-Third
666 International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 6642–6650.
667 International Joint Conferences on Artificial Intelligence Organization, 8 2024. doi:
10.24963/ijcai.2024/734. URL <https://doi.org/10.24963/ijcai.2024/734>. Main
Track.

668

669 Kaizhong Zheng, Shujian Yu, and Badong Chen. Ci-gnn: A granger causality-inspired
670 graph neural network for interpretable brain network-based psychiatric diagnosis. *Neural
671 Networks*, 172:106147, April 2024. doi: 10.1016/j.neunet.2024.106147. URL [https://
www.sciencedirect.com/science/article/abs/pii/S0893608024000716](https://www.sciencedirect.com/science/article/abs/pii/S0893608024000716).

672

673 Ao Zhou, Jianlei Yang, Yingjie Qi, Tong Qiao, Yumeng Shi, Cenlin Duan, Weisheng Zhao,
674 and Chunming Hu. Hgnas: Hardware-aware graph neural architecture search for edge
675 devices. *IEEE Transactions on Computers*, 73(12):2693–2707, 2024. doi: 10.1109/TC.
676 2024.3449108.

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701