

---

# Position: Generalist Medical AI Requires Community-Governed Infrastructure

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This position paper argues that Generalist Medical AI (GMAI) cannot be achieved  
2 through model scaling alone, because the key bottleneck is not model capability  
3 but the absence of shared infrastructure for coordinating distributed medical  
4 intelligence. High-value medical data, clinical expertise, validation signals, and  
5 deployment feedback are decentralized, privacy-sensitive, and institutionally gov-  
6 erned; this makes centralized approaches structurally insufficient. We argue that  
7 GMAI should be developed as open, community-governed, expert-endorsed in-  
8 frastructure rather than as isolated proprietary models or fragmented specialist  
9 systems. We propose four necessary infrastructure mechanisms: resource dis-  
10 covery, living expert validation, governed deployment feedback, and sustainable  
11 community governance. Concrete instantiations include Doctor’s Last Exam (a  
12 living expert-endorsed benchmark), PatientHub (shared longitudinal patient cases),  
13 Medical MCP (governed tool protocols), Agentic Hospital (deployment environ-  
14 ments), and AI-native Health Information Systems. This position reframes GMAI  
15 as a socio-technical infrastructure agenda for trustworthy, scalable, and equitable  
16 medical intelligence.

## 17 1 Introduction

18 Recent progress in foundation models [Bommasani et al., 2021] has encouraged a scaling-centered  
19 view of Generalist Medical AI [Moor et al., 2023]: *larger models, broader modalities, and more*  
20 *tasks* may eventually yield systems that support diagnosis, triage, imaging, documentation, and care  
21 coordination. This view is powerful but incomplete.

22 **Definition 1** *Generalist Medical AI is defined as a class of medical AI systems capable of generalizing*  
23 *across tasks, modalities, and clinical scenarios.*

24 **GMAI Challenges** Medicine differs from many general AI domains because its most valuable  
25 resources are structurally distributed. Medical data is highly sensitive, deeply professional, strongly  
26 regulated, and closely tied to patient privacy [Price and Cohen, 2019]. Much of the most valuable  
27 medical data remains distributed across hospital servers, research institutions, public health systems,  
28 and local clinical databases [Johnson et al., 2023]. These “sleeping data” and institutional silos make  
29 it unrealistic for any single company or centralized platform to collect, standardize, and govern all  
30 the data required for GMAI [Rieke et al., 2020]. Similarly, medical expertise is decentralized across  
31 specialties, institutions, regions, and clinical communities, making it difficult for proprietary systems  
32 to internalize all necessary knowledge in a scalable way. The development of GMAI therefore  
33 follows two complementary routes: 1) top-down scaling by large organizations and 2) bottom-up  
34 specialization by clinical communities. We argue that GMAI could be achieved through a *general-*  
35 *specialist integration* approach where generalist models provide broad multi-modal perception and  
36 high-level reasoning, while specialist components provide depth and precision [Nori et al., 2023b].



Figure 1: The challenge of sleeping medical data: high-value medical data and expertise are decentralized across hospitals, research institutions, public health systems, and clinical communities, making centralized collection unrealistic.

37 The key challenge is connecting these distributed components into a unified, trustworthy, and scalable  
 38 ecosystem.

39 **Our Position.** We argue that such integration requires a community-driven mechanism. The central  
 40 position of this paper is:

41 **GMAI will fail if treated as a model-scaling race; instead, it should be built as**  
 42 **open, community-governed, expert-endorsed infrastructure that coordinates**  
 43 **fragmented medical resources and expertise.**

44 Community-driven GMAI is both a technical necessity and a governance response. Technically,  
 45 trustworthy GMAI requires high-quality data, expert annotation, clinical feedback, rigorous bench-  
 46 marking [Jin et al., 2021, Pal et al., 2022, Wang et al., 2024a], continuous evaluation, evidence-  
 47 grounded retrieval [Lewis et al., 2020, Zakka et al., 2024], and human-in-the-loop deployment—all of  
 48 which depend on distributed expertise and community participation. From a governance perspective,  
 49 without open infrastructure, medical intelligence may become concentrated in proprietary systems  
 50 whose incentives may not align with global public health needs [Reddy et al., 2020, World Health  
 51 Organization, 2021]. An open, community-driven infrastructure can ensure that the benefits of  
 52 medical AI are shared broadly rather than limited to a few companies, countries, or privileged groups.

## 53 2 Related Work

54 **From Specialized to Generalist Medical AI** Early medical AI systems were typically built for  
 55 narrowly defined tasks, such as skin cancer classification [Esteva et al., 2017], chest X-ray diagno-  
 56 sis [Rajpurkar et al., 2017], and breast cancer screening [McKinney et al., 2020]. Despite strong  
 57 performance in controlled settings, each system often required task-specific data, models, workflows,  
 58 and deployment pipelines. Foundation models have shifted attention toward Generalist Medical  
 59 AI (GMAI), where one system can operate across tasks, modalities, and clinical contexts [Moor  
 60 et al., 2023]. Representative efforts include Med-PaLM and Med-PaLM 2 for medical question  
 61 answering and licensing exams [Singhal et al., 2023, 2025], Med-Gemini for multimodal medical  
 62 reasoning [Saab et al., 2024], and studies of general-purpose models such as GPT-4 in medical set-  
 63 tings [Nori et al., 2023a]. Open medical models, including HuatuoGPT, PMC-LLaMA, LLaVA-Med,  
 64 and BiomedGPT, further explore domain adaptation and multimodal understanding [Zhang et al.,  
 65 2023a, Wu et al., 2024, Li et al., 2023, Zhang et al., 2023b]. In contrast, we argue that GMAI is not  
 66 only a model-scaling problem, but also an infrastructure problem: it requires an open ecosystem that  
 67 connects general models, specialist tools, clinical knowledge bases, and local workflows.

68 **Medical Data, Knowledge, and Evaluation Infrastructure** A core barrier to medical AI is that  
 69 clinical data remains distributed across hospitals, research institutions, and public health systems.

70 Federated learning offers one path to use such data without centralization [Rieke et al., 2020, Dayan  
71 et al., 2021], while public datasets such as MIMIC-IV [Johnson et al., 2023] and resources such as  
72 UMLS [Bodenreider, 2004] support research but capture only part of real-world medicine. Retrieval-  
73 augmented generation links models to external evidence [Lewis et al., 2020], and systems such  
74 as Almanac suggest that guideline grounding can improve clinical reliability [Zakka et al., 2024,  
75 Xiong et al., 2024]. Evaluation has also advanced through benchmarks such as MedQA, PubMedQA,  
76 MedMCQA, CMB, and MedBench [Jin et al., 2021, 2019, Pal et al., 2022, Hendrycks et al., 2021,  
77 Wang et al., 2024a, Cai et al., 2024, Tu et al., 2024]. Yet many benchmarks remain static and  
78 exam-oriented, falling short on clinical judgment, safety, evidence use, and workflow fit. We therefore  
79 view data, knowledge, and evaluation as shared infrastructure that requires expert input, provenance  
80 tracking, and community maintenance.

81 **Medical Agents and Open Governance** Recent work explores medical AI in agentic and workflow-  
82 based settings. Agent Hospital simulates hospital roles and workflows for training medical agents [Li  
83 et al., 2024], AgentClinic evaluates agents in simulated clinical encounters [Schmidgall et al., 2024],  
84 and MedAgents explores multi-agent collaboration [Tang et al., 2024]. These works show that medical  
85 AI must interact with tools, records, and institutional workflows beyond text generation. Healthcare  
86 AI also raises governance questions around accountability, transparency, safety, and equity [Reddy  
87 et al., 2020, World Health Organization, 2021]. We extend these discussions by proposing open,  
88 community-driven GMAI infrastructure with shared licenses, expert-endorsed resources, community  
89 benchmarks, feedback channels, and sustainable institutions that align medical AI with public health  
90 needs rather than proprietary interests alone.

91 **Gap in Existing Work.** While these works advance medical AI capabilities, evaluation methods,  
92 and agentic systems, a critical gap remains: *no existing work proposes integrated infrastructure for*  
93 *coordinating distributed medical resources, expert validation, and governed deployment.* Federated  
94 learning addresses data privacy but not benchmarking, knowledge bases, or feedback loops. Open  
95 models provide technical artifacts but not governance mechanisms or quality assurance. Medical  
96 agents demonstrate workflow capabilities but lack real-world deployment infrastructure. Our frame-  
97 work fills this gap by proposing community-driven infrastructure that connects distributed data access,  
98 living expert validation, governed deployment feedback, and sustainable governance into a unified  
99 ecosystem for trustworthy GMAI.

## 100 3 Why GMAI needs Community-driven Infrastructure, Not Only Models

### 101 3.1 Why Top-down Scaling Alone is Insufficient

102 The development of medical AI has followed two complementary routes. The first is a top-down  
103 scaling route, where large organizations build powerful foundation models using advantages in  
104 data, compute, and infrastructure [Singhal et al., 2023, Saab et al., 2024, Nori et al., 2023a]. The  
105 second is a bottom-up specialization route, where researchers, clinicians, startups, and hospitals  
106 build models, datasets, tools, and workflows for specific diseases, specialties, departments, or care  
107 settings [Zhang et al., 2023a, Chen et al., 2023, Wu et al., 2024, Li et al., 2023]. In many general  
108 AI domains, top-down scaling can absorb large amounts of public or easily standardized data. In  
109 medicine, however, this assumption breaks down.

110 **Challenges of Top-down Scaling in GMAI** Top-down scaling assumes that large organizations  
111 can accumulate enough data, compute, and infrastructure to build broadly capable medical models.  
112 This assumption breaks down in medicine: high-value data is sensitive, regulated, institutionally  
113 governed, and often remains “sleeping” across hospitals, research institutions, public health systems,  
114 specialist centers, and regional healthcare networks [Price and Cohen, 2019, Johnson et al., 2023].  
115 Such data are hard to centralize due to privacy, incentives, incompatible formats, consent, and local  
116 governance constraints [Rieke et al., 2020, Dayan et al., 2021]. Expertise is likewise distributed  
117 across specialties, regions, and care roles. Centralized efforts such as Watson Health, and the gap  
118 between exam-benchmark performance and clinical reliability in frontier models, suggest that the  
119 bottleneck is not model capability alone, but distributed data access and expert validation.

120 **General-specialist Integration for GMAI.** The practical roadmap of GMAI is unlikely to be  
121 purely general or purely specialized. Generalist models can provide broad reasoning, multimodal  
122 perception, natural language interaction, and task coordination. Specialist models, disease-specific  
123 knowledge bases, clinical tools, local protocols, and expert systems provide depth, precision, and

Table 1: Comparison of approaches to medical AI development. Community-driven GMAI infrastructure could address limitations of both centralized proprietary systems and open models alone.

Aspect	Centralized Proprietary	Open Models Only	Community-driven GMAI
Data access	Limited to partnerships	Public datasets only	Federated + governed registries
Expertise	Internal teams	Volunteer contributors	Credentialed experts + incentives
Evaluation	Internal validation	Static benchmarks	Living expert-endorsed evaluation
Deployment	Controlled rollout	Uncontrolled release	Risk-stratified governed access
Governance	Corporate decisions	No formal governance	Community-endorsed mechanisms
Sustainability	Profit-driven	Volunteer-dependent	Reciprocal value creation

124 adaptation [Nori et al., 2023b]. The challenge is not whether generalist or specialist systems will  
 125 prevail, but how to connect them into a trustworthy, scalable, and governable infrastructure.

### 126 3.2 Reasons for Community-driven Infrastructure

127 General-specialist integration cannot be achieved by model composition alone. It requires a shared  
 128 infrastructure that coordinates generalist models, specialist resources, clinical tools, expert knowledge,  
 129 benchmarks, and deployment feedback. Such infrastructure is not merely a collection of open-source  
 130 models or volunteer contributions, but a set of governed mechanisms through which distributed  
 131 medical resources can be discovered, validated, connected, and improved over time.

132 Therefore, GMAI requires more than larger models, open releases, or isolated specialist systems. It  
 133 requires community-driven infrastructure with three core necessities. This infrastructure is also a  
 134 governance response: without open and accountable mechanisms, medical intelligence may become  
 135 concentrated in a few proprietary platforms whose incentives may not align with public health  
 136 needs [Reddy et al., 2020, World Health Organization, 2021]. Notably, even proprietary organiza-  
 137 tions have increasingly recognized the need for community-like expert participation, as shown by  
 138 HealthBench [Arora et al., 2025] and HealthBench Professional [Hicks et al., 2026], which draw on  
 139 the expertise of hundreds of physicians across multiple countries.

140 **Necessity I: Leveraging Segmented Resources** Without mechanisms for accessing distributed med-  
 141 ical resources, GMAI cannot move beyond narrow public datasets or proprietary hospital partnerships.  
 142 Medical data and expertise are segmented across hospitals, specialties, regions, public health systems,  
 143 and patient communities, and cannot be fully centralized because of privacy, consent, institutional  
 144 governance, and local clinical context. Community infrastructure should therefore provide metadata  
 145 registries, provenance tracking, privacy-preserving access, expert directories, and shared access rules,  
 146 so that locally governed resources become visible without being extracted into a single platform.

147 **Necessity II: Expert-endorsed Living Validation** Medical AI cannot be trusted through static  
 148 benchmarks or automatic metrics alone. Trustworthy GMAI requires expert annotation, clinical  
 149 judgment, safety review, evidence grounding, and continuous benchmark updates [Jin et al., 2021, Pal  
 150 et al., 2022, Wang et al., 2024a, Tu et al., 2024]. Community-driven evaluation turns benchmarking  
 151 from a one-time leaderboard into a living validation process, where expert users continuously identify  
 152 unsafe responses, missing evidence, workflow failures, and region-specific needs. This makes  
 153 evaluation part of system improvement rather than a final reporting step.

154 **Necessity III: Governed Deployment Feedback** GMAI must improve from real-world use, but  
 155 clinical deployment feedback cannot be collected through uncontrolled logging or purely commercial  
 156 extraction. Hospitals, doctors, patients, medical schools, and startups need reciprocal mechanisms:  
 157 contributors provide cases, annotations, tools, or workflow feedback, and in return receive better  
 158 models, safer assistants, reusable benchmarks, and shared infrastructure. Governed feedback channels,  
 159 responsible licenses, audit trails, and risk-aware access can transform deployment from a one-  
 160 way extraction process into an accountable learning loop aligned with safety, sustainability, and  
 161 equity [Reddy et al., 2020, World Health Organization, 2021].

### 162 3.3 Addressing Common Objections

163 This position paper addresses four common objections.

164 **Objection 1: “Proprietary systems are safer because medical AI is high-risk.”** Closed de-  
 165 velopment may provide control, but control is not the same as safety. Medical AI safety depends

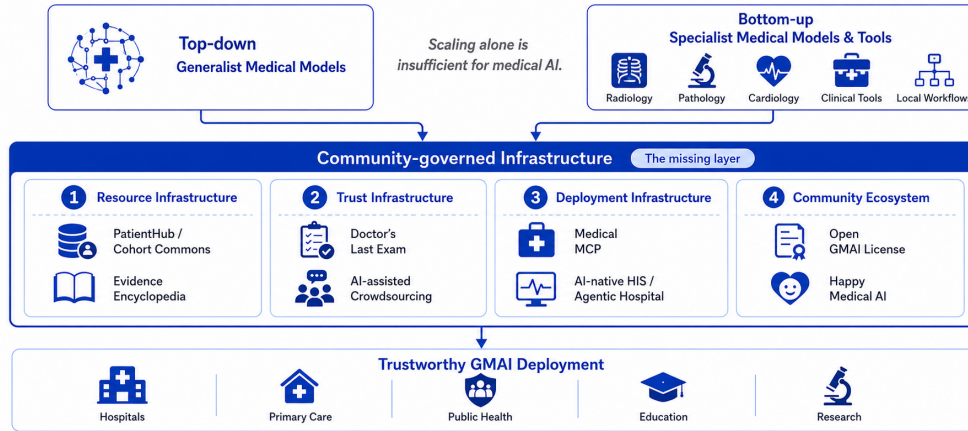


Figure 2: The Framework of GMAI.

166 on transparent evaluation, independent expert review, audit trails, evidence versioning, and post-  
 167 deployment monitoring [Reddy et al., 2020, World Health Organization, 2021]. These mechanisms  
 168 can be strengthened, not weakened, by community infrastructure. Openness should therefore mean  
 169 governed transparency, not uncontrolled release.

170 **Objection 2: “Federated learning already solves decentralized medical data.”** Federated  
 171 learning is important, but it solves only part of the problem [Rieke et al., 2020, Dayan et al., 2021].  
 172 GMAI also needs shared benchmarks, expert consensus, clinical knowledge bases, tool standards,  
 173 workflow environments, and continuous feedback. Federated learning might be a component of  
 174 community-driven GMAI, not a substitute for it.

175 **Objection 3: “Specialist models are enough; we do not need GMAI.”** Specialist models are  
 176 essential, but clinical practice is rarely confined to one specialty. Patients often have multiple diseases,  
 177 long histories, medications, social contexts, and operational constraints. The goal of GMAI is not to  
 178 replace specialists, but to coordinate them through general-specialist integration: broad reasoning  
 179 from generalist models combined with depth from specialist models, tools, and knowledge bases [Nori  
 180 et al., 2023b].

181 **Objection 4: “Open medical AI increases misuse risk.”** This concern is valid, but the answer  
 182 is not naive openness or complete closure. We advocate *governed openness*: credential-verified  
 183 contribution, risk-stratified access, responsible licensing, audit logging, and human approval for  
 184 high-risk actions. Such mechanisms make open medical AI more accountable while still allowing the  
 185 community to inspect, improve, and govern shared infrastructure.

## 186 4 Framework for Community-Driven GMAI

187 Community-driven Generalist Medical AI should not be understood as a single foundation model  
 188 released to the public. It is a shared infrastructure for turning distributed medical data, expert  
 189 knowledge, evaluation signals, deployment feedback, and governance capacity into reusable medical  
 190 intelligence. We organize this infrastructure into four interdependent pillars: (1) **Resource Infras-**  
 191 **tructure** for models, data, cohort resources, and expert knowledge; (2) **Trust Infrastructure**  
 192 for benchmarks, crowdsourced annotation, evidence validation, and feedback loops; (3) **Deployment**  
 193 **Infrastructure** for tools, agentic environments, agent platforms, and AI-native healthcare systems;  
 194 and (4) **Community Ecosystem** for education, entrepreneurship, governance, multilingual participa-  
 195 tion, and long-term maintenance. These pillars jointly provide a practical route toward trustworthy,  
 196 scalable, and globally accessible GMAI.

### 197 4.1 Resource Infrastructure: Models and Data

198 The first pillar provides foundational resources: open models, shared data, and expert knowledge.  
 199 These components are tightly coupled: models require data and knowledge, while data and knowledge  
 200 require models to become useful.

#### 201 4.1.1 Open-sourced Models

202 Open-sourced models [Wu et al., 2024, Zhang et al., 2023a, Chen et al., 2023, Touvron et al., 2023]  
203 are necessary but insufficient. Medical AI involves patient privacy, institutional consent, clinical  
204 safety, and deployment accountability [World Health Organization, 2021, Reddy et al., 2020]. A  
205 community-driven ecosystem requires governance specifying what can be shared, by whom, under  
206 what conditions, and with what obligations. Moreover, open models must support deployment across  
207 heterogeneous hardware ecosystems including NVIDIA, AMD or Ascend accelerators; so that GMAI  
208 infrastructure is not locked to a single vendor and can be deployed in regions with different hardware  
209 availability and regulatory requirements.

210 **Open GMAI License** A community-developed Open GMAI License can define how datasets,  
211 models, tools, benchmarks, and knowledge bases are shared, reused, and deployed. Community  
212 funding can support open-source organizations, annotation fellowships, hospital pilots, and global  
213 programs, transforming “sleeping medical data” into public-interest infrastructure.

#### 214 4.1.2 Shared Data

215 Rather than centralizing sensitive data, the community can build a registry where institutions publish  
216 metadata, provenance, access rules, and consent conditions. Actual data remains locally stored  
217 or accessed through privacy-preserving computation, while the registry records what exists, who  
218 contributed it, and how it may be used.

219 **PatientHub: Longitudinal Cases and Cohort Commons.** A representative example is PatientHub,  
220 a governed infrastructure for longitudinal patient cases and cohort-level resources. Each case  
221 records a patient’s care trajectory across symptoms, visits, diagnoses, laboratory tests, imaging,  
222 medications, interventions, outcomes, and follow-up. These cases can be organized into cohorts by  
223 disease, phenotype, treatment pathway, care setting, demographic profile, or outcome, supporting  
224 model training, evaluation, medical education, comparative analysis, and workflow simulation.  
225 PatientHub therefore turns individual care trajectories into reusable cohort resources while preserving  
226 privacy through de-identification, synthetic or hybrid case construction, provenance tracking, consent  
227 documentation, and clear intended-use rules.

#### 228 4.2 Trust Infrastructure: Benchmarking and Crowdsourcing

229 Trust infrastructure ensures GMAI resources meet clinical standards through rigorous evaluation and  
230 expert validation, combining benchmarking and crowdsourced annotation.

##### 231 4.2.1 Benchmark

232 Medical evaluation requires expert judgment and real clinical context beyond static benchmarks or  
233 automatic metrics. A trustworthy infrastructure combines benchmark construction, expert review,  
234 post-deployment feedback, and continuous updates.

235 **Doctor’s Last Exam.** Inspired by expert-level benchmarks designed to evaluate frontier AI capa-  
236 bilities across broad academic domains, a crowd-funded “Doctor’s Last Exam” becomes a living  
237 examination program jointly maintained by clinicians, medical schools, and hospitals. Cases cover  
238 diagnosis, triage, longitudinal care, multimodal interpretation, and workflow execution. Doctors  
239 correct unsafe responses, patients report usability issues, and hospitals provide workflow feedback,  
240 creating a loop from real-world needs to benchmark design and system improvement.

##### 241 4.2.2 Crowdsourcing

242 Medical annotations require clinical training, specialty knowledge, and awareness of uncertainty. A  
243 platform should verify credentials, route tasks to suitable experts, and record annotator background  
244 and confidence. It supports dataset construction, benchmark authoring, RAG validation, and feedback  
245 collection.

246 **AI-Assisted Crowdsourcing Platform.** AI systems draft labels, extract findings, and propose  
247 diagnoses. Verified experts review, correct, or approve outputs, focusing on high-value judgment  
248 rather than repetitive screening. The platform supports audit trails, privacy review, disagreement

249 resolution, and quality control. This creates a loop where AI prepares annotations, experts provide  
250 judgment, corrected labels improve models, and improved systems assist future annotation.

251 **AI-native Evidence-based Encyclopedia-like Guidelines.** An AI-native medical encyclopedia  
252 should be community-owned, organized by specialty, and governed by expert-endorsed editorial  
253 structures. Each claim carries provenance, evidence level, update history, and clinical scope. The  
254 encyclopedia supports retrieval, citation, version control, and API access, enabling GMAI to rely  
255 on current evidence rather than static memory. Research access remains free, while contributor  
256 incentives tie to usage records. Revenue from clinical or educational use returns to authors and the  
257 community.

### 258 **4.3 Deployment Infrastructure: Agentic Tools and Environments**

259 The third pillar enables GMAI to move from research prototypes to deployable systems. This layer  
260 provides the technical infrastructure for medical agents to interact with tools, hospital systems, and  
261 real-world workflows.

#### 262 **4.3.1 Medical Tool (MCP)**

263 GMAI requires a shared medical tool layer where community-validated APIs, calculators, and  
264 workflow interfaces can be safely used by agents. Models should call tools for current evidence,  
265 calculations, safety checks, or human approval. Standardized protocols enable tools from hospitals,  
266 specialty groups, and startups to be discovered, invoked, and audited under common rules.

267 **Medical MCP** A community-based medical API hub documents each tool with its schema, provenance,  
268 validation status, risk level, and required permission. Low-risk tools may be called automati-  
269 cally; high-risk tools require explicit confirmation or expert review. This allows agents to combine  
270 general reasoning with trustworthy utilities while making invocations traceable and governable.

#### 271 **4.3.2 Agentic Environments**

272 GMAI should be deployed as agentic systems capable of routing tasks, invoking tools, retrieving  
273 knowledge, and coordinating workflows [Li et al., 2024, Schmidgall et al., 2024]. This requires  
274 environments with shared tools, standardized patient simulation, and HIS interfaces to train, evaluate,  
275 and safely deploy agents before clinical use.

276 **Agentic Hospital** An agentic hospital environment simulates hospital roles, departments, patient  
277 trajectories, devices, and HIS interactions. It supports training, evaluation, safety testing, and  
278 workflow design through reusable standardized patient simulators. At deployment, this layer connects  
279 agents to AI-native HIS infrastructure with secure, permission-controlled, auditable interfaces and  
280 routing mechanisms [Tang et al., 2024].

281 **AI-native HIS** AI-native Health Information Systems (HIS) differs from Agentic Hospital en-  
282 vironments. Agentic Hospital provides simulated organizational settings for training and testing;  
283 AI-native HIS provides agent-compatible software infrastructure for interacting with healthcare data  
284 and workflows. Unlike traditional systems designed for human interfaces, AI-native HIS supports  
285 structured APIs, tool-calling protocols, access control, audit logging, and human approval. Together,  
286 they create the operating environment for testable, deployable agentic GMAI.

287 **Agent Platform** A community GMAI portal allows doctors to test copilots, hospitals to evaluate  
288 workflow agents, students to train with standardized patients, and researchers to access benchmarks.  
289 It collects real demands: where models fail, what evidence is missing, which workflows are hard to  
290 automate. For willing contributors, it provides governed channels for sharing de-identified cases,  
291 feedback logs, and annotations under explicit consent and privacy review. The platform becomes  
292 both an access point and a feedback engine for improving open medical intelligence.

### 293 **4.4 Community Ecosystem: Education, Entrepreneurship, and Governance**

294 The fourth pillar ensures long-term viability of the GMAI ecosystem. This addresses community  
295 building, economic sustainability, and institutional governance that transforms GMAI into a durable  
296 public infrastructure.

#### 297 4.4.1 Medical Education

298 GMAI requires an open medical AI education ecosystem. Inspired by communities like Datawhale,  
299 the community can build a “Happy Medical AI” initiative: a free, practice-oriented curriculum for  
300 doctors, students, patients, developers, and public health workers to understand, use, evaluate, and  
301 contribute to medical AI. This layer combines online courses, tutorials, benchmark challenges, and  
302 contribution portals with offline bootcamps, hospital workshops, annotation camps, and hackathons  
303 to train contributors, collect needs, validate tools, and build trust.

304 **Online Tutorial (Happy Medical AI).** Inspired by community-driven open education projects,  
305 Happy Medical AI provides modular courses, clinical case exercises, annotation practice, benchmark  
306 tasks, agent-building examples, and deployment playbooks. Different tracks serve different com-  
307 munities: clinicians learn evidence retrieval and workflow support; students practice reasoning with  
308 standardized patients; developers learn data governance, tool APIs, and evaluation; patients learn  
309 health literacy and safe AI use. Education becomes an entry point for participation, turning learners  
310 into contributors.

#### 311 4.4.2 Open Entrepreneurship and Community Governance

312 Beyond research and technical infrastructure, GMAI requires sustainable governance mechanisms.

313 **Open Entrepreneurship.** Dedicated medical AI seed funds and community-based startup programs  
314 can help early-stage teams build on shared GMAI infrastructure. Potential directions include specialty  
315 agents, AI-native HIS tools, patient-facing applications, education platforms, workflow automation,  
316 multilingual services, annotation platforms, and privacy-preserving deployment. The goal is to  
317 create a healthy ecosystem where open infrastructure enables responsible innovation while startups  
318 contribute improvements, tools, and feedback back to the commons.

319 **Global Community and Multilingual Support.** A truly community-driven GMAI must be glob-  
320 ally inclusive. This requires multilingual model support [Wang et al., 2024b, Zheng et al., 2024],  
321 localized medical guidelines, cross-cultural clinical knowledge, and regional community chapters.  
322 The infrastructure should support major languages (English, Chinese, Spanish, Arabic, Hindi, etc.)  
323 and enable local communities to contribute region-specific medical knowledge, traditional medicine  
324 practices, and culturally appropriate healthcare workflows. Global coordination mechanisms (includ-  
325 ing regional hubs, international workshops, and cross-border collaboration platforms) can connect  
326 medical professionals, researchers, and developers worldwide, ensuring that GMAI serves diverse  
327 populations rather than replicating existing healthcare inequalities.

## 328 5 Call for Actions in GMAI

### 329 5.1 What We Should Do Now

330 Community-driven GMAI will not emerge spontaneously. It requires deliberate action from specific  
331 stakeholders. We outline concrete next steps.

332 **ML researchers** should prioritize building shared evaluation infrastructure over publishing yet  
333 another medical benchmark paper. Contribute cases to Doctor’s Last Exam, develop evaluation  
334 protocols that test clinical safety and workflow fitness (not just accuracy on multiple-choice questions),  
335 and release evaluation code and data under open licenses. When developing medical AI models,  
336 support deployment across heterogeneous hardware—including both NVIDIA GPUs and domestic  
337 accelerators such as Huawei Ascend—to avoid vendor lock-in.

338 **Hospitals and health systems** should participate in data registries that publish metadata, provenance,  
339 and access conditions without centralizing sensitive data. Provide deployment feedback on where AI  
340 systems fail in real workflows. Contribute de-identified cases to shared repositories under governed  
341 consent frameworks. The return is access to better models, benchmarks, and tools built on the  
342 collective infrastructure.

343 **Industry** should adopt open licensing frameworks (Open GMAI License) that define responsible  
344 sharing, reuse, and deployment obligations. Contribute tools to Medical MCP with documented  
345 schemas and validation status. Fund open infrastructure as shared pre-competitive investment rather  
346 than treating all medical AI as proprietary advantage.

347 **Medical schools and professional societies** should integrate AI literacy into medical education, con-  
348 tribute expert annotations through credentialed platforms, and participate in benchmark governance.  
349 Medical expertise is the scarcest resource in GMAI development; professional communities must  
350 organize to provide it sustainably.

351 **Policymakers** should develop regulatory frameworks that enable cross-border medical data gover-  
352 nance, recognize community-endorsed evaluation as a pathway to clinical validation, and fund open  
353 medical AI infrastructure as public health investment. Current regulations often assume centralized  
354 development; policy must evolve to support distributed, community-driven approaches.

## 355 5.2 Applications

356 GMAI infrastructure enables diverse applications across healthcare settings. In hospitals, it supports  
357 clinical workflows, imaging interpretation, and evidence-grounded decision support [Topol, 2019, Tu  
358 et al., 2024]. For patients, it enables continuous longitudinal care, chronic disease management, and  
359 access to medical intelligence in underserved regions [Thirunavukarasu et al., 2023]. In public health,  
360 it facilitates population-scale disease surveillance and emergency response [Dayan et al., 2021, Rieke  
361 et al., 2020]. For education and research, it provides scalable learning infrastructure and accelerates  
362 biomedical innovation [Schmidgall et al., 2024, Jin et al., 2021]. Finally, it can serve as the reasoning  
363 layer for embodied healthcare agents including surgical and rehabilitation robots [Li et al., 2024].  
364 This broad application landscape demonstrates why GMAI must be community-driven: no single  
365 institution can cover all healthcare scenarios, specialties, workflows, and populations.

## 366 6 Challenges and Limitations

367 We acknowledge significant challenges as below that must be addressed for GMAI vision to succeed.

368 **Scope of This Work.** As a position paper, our proposed framework and programs are forward-  
369 looking and have not yet been fully implemented or empirically validated. The concrete programs  
370 (e.g., Doctor’s Last Exam and PatientHub) represent aspirational designs whose feasibility depends  
371 on sustained community participation, institutional buy-in, and adequate funding. We present this  
372 framework as a research agenda and call to action rather than a validated solution.

373 **Coordination Costs and Governance Complexity.** Open communities face inherent coordination  
374 overhead. Aligning diverse stakeholders (including hospitals, researchers, clinicians, patients, regula-  
375 tors, and industry) across different countries, legal systems, and incentive structures is substantially  
376 harder than centralized development. Decision-making may be slow, and governance structures may  
377 struggle to balance inclusivity with efficiency.

378 **Quality Control and Safety Risks.** Open contributions introduce quality variance. Unlike prop-  
379 rietary systems with unified quality pipelines, community-contributed data, annotations, and tools  
380 may vary in reliability. In healthcare, quality failures carry patient safety risks. Our framework  
381 addresses this through expert endorsement, credential verification, and tiered access controls, but  
382 these mechanisms add friction and may not catch all errors.

383 **Privacy and Regulatory Heterogeneity.** Medical data governance varies dramatically across  
384 jurisdictions. GDPR, HIPAA, China’s Personal Information Protection Law, and other frameworks  
385 impose different requirements on data sharing, consent, and cross-border transfer. Building truly  
386 global infrastructure while respecting all regulatory regimes is an unsolved challenge that may limit  
387 the scope of data sharing in practice.

## 388 7 Conclusion

389 Generalist Medical AI should not be pursued through model scaling alone. Because medical  
390 data, expertise, trust, and deployment feedback are inherently distributed, GMAI requires an open,  
391 community-driven, and expert-endorsed infrastructure. This paper argues for building such infrastruc-  
392 ture around shared resources, trustworthy evaluation, agentic deployment, and sustainable governance.  
393 By connecting generalist models with specialist knowledge, clinical feedback, responsible licens-  
394 ing, and community participation, community-driven GMAI offers a practical path toward scalable,  
395 trustworthy, and equitable medical intelligence for global healthcare.

## 396 References

- 397 Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela,  
398 Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, et al. Health-  
399 bench: Evaluating large language models towards improved human health. *arXiv preprint*  
400 *arXiv:2505.08775*, 2025.
- 401 Olivier Bodenreider. The Unified Medical Language System (UMLS): integrating biomedical  
402 terminology. *Nucleic Acids Research*, 32(suppl\_1):D267–D270, 2004.
- 403 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arber, Sydney von Arx, et al.  
404 On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 405 Yan Cai, Linlin Wang, Ye Wang, Gerard de Melo, Ya Zhang, Yanfeng Wang, and Liang He. Med-  
406 Bench: A large-scale Chinese benchmark for evaluating medical large language models. In *AAAI*,  
407 2024.
- 408 Junying Chen, Xidong Wang, Anningzhe Gao, Feng Jiang, Shunian Chen, Hongbo Zhang,  
409 Dingjie Song, Wenya Xie, Chuyi Kong, Jianquan Li, Xiang Wan, Haizhou Li, and Benyou  
410 Wang. HuatuoGPT-II, one-stage training for medical adaption of LLMs. In *arXiv preprint*  
411 *arXiv:2311.09774*, 2023.
- 412 Ittai Dayan, Holger R Roth, Aoxiao Zhong, Ahmed Harouni, Amilcare Gentili, Anas Z Abidin, et al.  
413 Federated learning for predicting clinical outcomes in patients with COVID-19. *Nature Medicine*,  
414 27(10):1735–1743, 2021.
- 415 Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and  
416 Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks.  
417 *Nature*, 542(7639):115–118, 2017.
- 418 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
419 Steinhardt. Measuring massive multitask language understanding. In *ICLR*, 2021.
- 420 Rebecca Soskin Hicks, Mikhail Trofimov, Dominick Lim, Rahul K Arora, Foivos Tsimpourlas,  
421 Preston Bowman, Michael Sharman, Chi Tong, Kavin Karthik, Arnav Dugar, et al. Health-  
422 bench professional: Evaluating large language models on real clinician chats. *arXiv preprint*  
423 *arXiv:2604.27470*, 2026.
- 424 Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. What  
425 disease does this patient have? A large-scale open domain question answering dataset from medical  
426 exams. *Applied Sciences*, 11(14):6421, 2021.
- 427 Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. PubMedQA: A  
428 dataset for biomedical research question answering. In *EMNLP*, 2019.
- 429 Alistair E W Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng,  
430 Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. MIMIC-IV, a freely accessible  
431 electronic health record dataset. *Scientific Data*, 10(1):1, 2023.
- 432 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
433 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-  
434 tion for knowledge-intensive NLP tasks. In *NeurIPS*, 2020.
- 435 Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan  
436 Naumann, Hoifung Poon, and Jianfeng Gao. LLaVA-Med: Training a large language-and-vision  
437 assistant for biomedicine in one day. In *NeurIPS*, 2023.
- 438 Junkai Li, Yunghwei Lai, Weitao Li, Jingyi Ren, Meng Zhang, Xinhui Kang, Siyu Wang, Peng Li,  
439 Ya-Qin Zhang, Weizhi Ma, and Yang Liu. Agent hospital: A simulacrum of hospital with evolvable  
440 medical agents. *arXiv preprint arXiv:2405.02957*, 2024.
- 441 Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova,  
442 Hutan Ashrafian, et al. International evaluation of an AI system for breast cancer screening. *Nature*,  
443 577(7788):89–94, 2020.

- 444 Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec,  
445 Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence.  
446 *Nature*, 616(7956):259–265, 2023.
- 447 Harsha Nori, Nicholas King, Scott Mayer McKinney, Dean Carignan, and Eric Horvitz. Capabilities  
448 of GPT-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*, 2023a.
- 449 Harsha Nori, Yin Tat Lee, Sheng Zhang, Dean Carignan, Richard Edgar, Nicolo Fusi, Nicholas King,  
450 Jonathan Larson, Yuanzhi Li, Weishung Liu, et al. Can generalist foundation models outcompete  
451 special-purpose tuning? Case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023b.
- 452 Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. MedMCQA: A large-scale  
453 multi-subject multi-choice dataset for medical domain question answering. In *Conference on*  
454 *Health, Inference, and Learning (CHIL)*, 2022.
- 455 W Nicholson Price and I Glenn Cohen. Privacy in the age of medical big data. *Nature Medicine*, 25  
456 (1):37–43, 2019.
- 457 Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy  
458 Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, Matthew P Lungren, and Andrew Y Ng.  
459 CheXNet: Radiologist-level pneumonia detection on chest X-rays with deep learning. *arXiv*  
460 *preprint arXiv:1711.05225*, 2017.
- 461 Sandeep Reddy, Sonia Allan, Simon Coghlan, and Paul Cooper. A governance model for the  
462 application of AI in health care. *Journal of the American Medical Informatics Association*, 27(3):  
463 491–497, 2020.
- 464 Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon  
465 Bakas, et al. The future of digital health with federated learning. *npj Digital Medicine*, 3:119,  
466 2020.
- 467 Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, et al. Capabili-  
468 ties of Gemini models in medicine. *arXiv preprint arXiv:2404.18416*, 2024.
- 469 Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor.  
470 AgentClinic: A multimodal agent benchmark to evaluate AI in simulated clinical environments.  
471 *arXiv preprint arXiv:2405.07960*, 2024.
- 472 Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan  
473 Scales, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180,  
474 2023.
- 475 Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Mohamed Amin, Le Hou, et al.  
476 Toward expert-level medical question answering with large language models. *Nature Medicine*, 31  
477 (3):943–950, 2025.
- 478 Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman Cohan,  
479 and Mark Gerstein. MedAgents: Large language models as collaborators for zero-shot medical  
480 reasoning. In *Findings of ACL*, 2024.
- 481 Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang  
482 Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*, 29(8):  
483 1930–1940, 2023.
- 484 Eric J Topol. High-performance medicine: the convergence of human and artificial intelligence.  
485 *Nature Medicine*, 25(1):44–56, 2019.
- 486 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
487 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLaMA: Open and  
488 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 489 Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang,  
490 et al. Towards generalist biomedical AI. *NEJM AI*, 1(3), 2024.

- 491 Xidong Wang, Guiming Hardy Chen, Dingjie Song, Zhiyi Zhang, Zhihong Chen, Qingying Xiao,  
492 Feng Jiang, Jianquan Li, Xiang Wan, Benyou Wang, and Haizhou Li. CMB: A comprehensive  
493 medical benchmark in Chinese. In *NAACL*, 2024a.
- 494 Xidong Wang, Nuo Chen, Junyin Chen, Yidong Wang, Guorui Zhen, Chunxian Zhang, Xiangbo  
495 Wu, Yan Hu, Anningzhe Gao, Xiang Wan, et al. Apollo: A lightweight multilingual medical llm  
496 towards democratizing medical ai to 6b people. *arXiv preprint arXiv:2403.03640*, 2024b.
- 497 World Health Organization. Ethics and governance of artificial intelligence for health: WHO guidance.  
498 Technical report, World Health Organization, 2021.
- 499 Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. PMC-LLaMA:  
500 Toward building open-source language models for medicine. *Journal of the American Medical  
501 Informatics Association*, 31(9):1833–1843, 2024.
- 502 Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. Benchmarking retrieval-augmented  
503 generation for medicine. In *Findings of ACL*, 2024.
- 504 Cyril Zakka, Rohan Shad, Akash Chaurasia, Alex R Dalal, Jennifer L Kim, Michael Moor, Robyn  
505 Fong, Curran Phillips, Kevin Alexander, Euan Ashley, et al. Almanac—retrieval-augmented  
506 language models for clinical medicine. *NEJM AI*, 1(2), 2024.
- 507 Hongbo Zhang, Junying Chen, Feng Jiang, Fei Yu, Zhihong Chen, Guiming Chen, Jianquan Li, Xi-  
508 angbo Wu, Zhiyi Zhang, Qingying Xiao, Xiang Wan, Benyou Wang, and Haizhou Li. HuatuoGPT,  
509 towards taming language model to be a doctor. In *Findings of EMNLP*, 2023a.
- 510 Kai Zhang, Rong Zhou, Eashan Adhikarla, Zhiling Yan, Yixin Liu, Jun Yu, Xun Chen, et al.  
511 BiomedGPT: A unified and generalist biomedical generative pre-trained transformer for vision,  
512 language, and multimodal tasks. *arXiv preprint arXiv:2305.17100*, 2023b.
- 513 Guorui Zheng, Xidong Wang, Juhao Liang, Nuo Chen, Yuping Zheng, and Benyou Wang. Efficiently  
514 democratizing medical llms for 50 languages via a mixture of language family experts. *arXiv  
515 preprint arXiv:2410.10626*, 2024.