

# AI System Development

## Proposal: Hospital AI Assistant

### for Clinical Decision Support

**Student Name:** [Your Name]

**Student ID:** [Your ID]

**Course:** AI Governance and Ethics

**Professor:** [Professor Name]

**Date Submitted:** November 21, 2025

---

## Introduction

---

### *System Overview*

The proposed Hospital AI Assistant is a clinical decision support system designed to augment healthcare delivery through intelligent symptom triage, differential diagnosis suggestion, and post-discharge patient management. The system processes patient-reported symptom narratives, clinical histories, vital signs from wearable devices, and preliminary clinical data (imaging metadata, laboratory results pending physician interpretation) to provide real-time clinical recommendations to hospital staff, including physicians, nurses, and allied health professionals, as well as patient-facing guidance on symptom severity and appropriate care pathways.

### *System Purpose and Clinical Value*

Hospital emergency departments and acute care settings face systemic challenges: prolonged patient wait times, variable triage accuracy, delayed identification of high-risk patients, and occasional missed diagnoses due to cognitive load and information fragmentation. The Hospital AI Assistant addresses these gaps by:

1. **Improving triage efficiency:** Enabling faster, data-driven patient prioritization
2. **Supporting clinical decision-making:** Offering evidence-based differential diagnoses for physician review and validation
3. **Enhancing early detection:** Identifying patients at elevated risk for adverse outcomes through continuous monitoring
4. **Reducing diagnostic errors:** Providing systematic symptom analysis to complement clinical judgment

Critically, the system operates as an *assistive technology*, not an autonomous diagnostic tool. All high-risk clinical decisions require human physician confirmation before implementation.

## *User Categories and Use Cases*

### **Direct Users:**

- **Physicians:** Receive AI-generated clinical prompts, confidence scores, and recommended diagnostic pathways during in-clinic consultations
- **Nursing and Allied Health Staff:** Utilize symptom triage modules for initial patient intake, test prioritization, and post-discharge risk assessment
- **Patients:** Access smartphone interface displaying symptom severity alerts, evidence-based self-triage guidance, and post-discharge adherence reminders

### **System Inputs:**

- Patient-reported symptom descriptions (free text and structured symptom checklist)
- Medical history and comorbidities
- Vital sign data from bedside monitors and wearable devices
- Preliminary clinical metadata (pending imaging/laboratory interpretation, not final diagnostic reads)

# Potential Harms

---

## *Research on Similar AI Healthcare Systems*

Recent implementations of AI diagnostic assistants in hospital settings have revealed critical risks. The Tsinghua University AI Hospital pilot and emerging literature on AI-assisted clinical triage systems document recurring failure modes: representation bias in training data, overreliance on model outputs by clinical staff, and security vulnerabilities in health data handling. These real-world deployments provide essential evidence for structured harm analysis.

### *Identified Potential Harms*

#### **Harm 1: Privacy and Data Governance Failure**

**Description:** Large-scale collection, storage, and processing of identifiable health information (Protected Health Information under HIPAA/PIPEDA standards) without strict minimization and access controls creates vulnerability to data leakage, unauthorized access, and patient privacy violations.

**Stakeholders Affected:** Patients (data subjects), healthcare organizations, regulatory bodies

#### **Evidence and Context:**

- Microsoft Responsible AI Standard (2022a) designates medical AI systems as requiring comprehensive data protection impact assessments, data classification protocols, and role-based access control implementation during the design phase, specifically in the Accountability section (Goal A5)
- The Microsoft RAI Impact Assessment Guide (2022b) documents privacy as a foundational control for health data systems, requiring explicit governance frameworks that address data minimization, access control, and retention policies
- Healthcare sector research demonstrates that API-based clinical systems expand the attack surface for sensitive data, necessitating encryption, key management, and continuous access auditing

#### **Specific Risks:**

- Unauthorized access or breaches exposing patient medical histories and identities
- Retention of patient data beyond necessary operational periods

- Inadequate consent mechanisms or patient withdrawal options for data usage
- 

## Harm 2: Diagnostic Error and Clinical Safety Compromise

**Description:** Uncertainty in AI-generated recommendations, particularly for underrepresented patient populations, can result in misdiagnosis, under-triage of critical cases, or over-triggering of unnecessary interventions, directly impacting patient safety and clinical outcomes.

**Stakeholders Affected:** Patients, physicians, nursing staff, hospitals

### Evidence and Context:

- Microsoft Responsible AI Standard (2022a) addresses human oversight in Goal A5, which mandates that stakeholders responsible for managing and controlling the system must have capabilities to understand when and how to override, intervene, or interrupt the system. This requirement reflects the recognition that high-consequence clinical decisions cannot be fully automated
- Microsoft RAI standards require specific accuracy evaluation across demographic and clinical subgroups, with documented performance thresholds for deployment in safety-critical domains (Microsoft, 2022b)
- Recent AI hospital deployments position decision support systems explicitly as tools that augment rather than replace clinical judgment, acknowledging that human accountability in diagnosis is irreducible (Med-Tech World, 2023)

### Specific Risks:

- Sparse training data for minority patient populations, rare disease presentations, or clinical extremes (pediatric, geriatric cases) reduces model confidence and accuracy for vulnerable groups
  - High-symptom-variability presentations (e.g., myocardial infarction presenting atypically in women, or sepsis in immunocompromised patients) exceed model training distribution
  - False negatives in high-risk pathways (stroke, cardiac events, severe infection) create life-threatening delays in care
- 

## Harm 3: Bias and Health Equity Inequity

**Description:** Systemic underrepresentation of minority, low-income, and geographically remote populations in training datasets creates disparate model performance, amplifying existing healthcare disparities and perpetuating structural inequities.

**Stakeholders Affected:** Minority patient populations, low-income communities, rural/remote patients, women, pediatric and geriatric populations

#### Evidence and Context:

- Microsoft Responsible AI Standard (2022a) mandates fairness evaluation across identified demographic groups through Goals F1-F3, requiring documented performance gaps and mitigation strategies prior to deployment. These standards recognize that healthcare AI systems present particular risk to vulnerable populations due to downstream consequences of misclassification
- The RAI Impact Assessment Guide (2022b) requires explicit demographic fairness testing as a baseline compliance control, specifically noting that performance disparities in healthcare systems warrant conservative deployment strategies and continuous monitoring
- Healthcare AI literature emphasizes that demographic imbalance in training data—particularly underrepresentation of women, ethnic minorities, children, elderly, rural populations, and non-English speakers—creates documented performance degradation for these groups

#### Specific Risks:

- Training data imbalance: underrepresented groups have sparse or absent representation in historical clinical datasets
- Performance degradation: model accuracy, sensitivity, and specificity vary significantly across demographic groups, leading to under-triage or inappropriate intervention for vulnerable populations
- Compounding disadvantage: patients from socioeconomically disadvantaged backgrounds may have atypical presentation patterns (e.g., delayed healthcare-seeking, multiple comorbidities) that reduce model performance precisely where equity is most critical

---

## Harm 4: Over-Reliance and Blurred Accountability

**Description:** Clinical staff may unconsciously treat AI outputs as authoritative clinical guidelines, eroding independent clinical reasoning. Conversely, patients may treat triage recommendations as

definitive diagnoses, delaying necessary clinical evaluation or self-managing inappropriately.

**Stakeholders Affected:** Physicians, nursing staff, patients, hospital liability

#### Evidence and Context:

- RAI Impact Assessment Guide (2022b) requires explicit responsibility chains and human confirmation points in clinical decision workflows to prevent automation bias, which is defined as users' tendency to favor automated recommendations over their own judgment, particularly under time pressure
- Recent AI hospital deployments position systems strictly as support tools, with clear demarcation of human vs. machine accountability in clinical workflows (Tsinghua University, 2023; Med-Tech World, 2023)
- Human factors research documents that automation bias increases under conditions of high cognitive load, time pressure, and complexity—all factors present in emergency department settings

#### Specific Risks:

- Physicians may inadvertently defer clinical judgment to model outputs without adequate critical appraisal, weakening diagnostic reasoning
- Patients may delay seeking clinical evaluation if triage system indicates "low risk," even in the presence of evolving or serious symptoms
- Liability exposure: ambiguity regarding responsibility (provider vs. system vs. organization) for adverse outcomes following AI recommendations

---

## Harm 5: Security and Attack Surface Vulnerabilities

**Description:** AI models and clinical APIs are subject to adversarial attacks, model extraction, data poisoning, and ransomware targeting sensitive health information infrastructure, disrupting clinical operations and compromising system integrity.

**Stakeholders Affected:** Hospital IT infrastructure, patients, clinical operations, data security

#### Evidence and Context:

- Microsoft Responsible AI Standard (2022a) designates security, reliability, and anti-tampering as foundational requirements in Goal RS1-RS3, requiring anomaly detection, rate

limiting, and rapid incident response protocols for health systems handling protected data

- Healthcare cybersecurity research documents an increase in ransomware and targeted cyber attacks on medical records systems and clinical decision support infrastructure, with attacks specifically targeting API interfaces of AI systems (National Institute of Standards and Technology, 2024)
- API-based model deployment increases exposure to model stealing attacks (where adversaries reverse-engineer models through repeated queries), adversarial input attacks (where crafted inputs trigger unsafe outputs), and distributed denial-of-service attacks that disable clinical support when most needed

### **Specific Risks:**

- Model extraction attacks: adversaries reverse-engineer the model through repeated API queries
  - Adversarial inputs: carefully crafted symptom inputs trigger erratic or unsafe recommendations
  - Infrastructure compromise: ransomware or system outages disable clinical decision support when most needed
  - Data exfiltration: patient medical records and identifiable health information stolen or accessed without authorization
- 

## **Risk Mitigation**

---

### *Mitigation Framework*

Each identified harm is addressed through **technical, organizational, and societal/compliance** controls, with designated stakeholder owners and measurable success metrics. This layered approach aligns with Microsoft RAI Standard requirements (Goals A5, T1-T3, F1-F3, RS1-RS3) and reflects best practices in healthcare AI governance.

---

### *Mitigation 1: Privacy and Data Governance*

#### **Technical Controls**

## **TECHNICAL CONTROLS**

- **Data Minimization:** Collect only data fields clinically necessary for symptom triage and differential diagnosis; explicit technical architecture to exclude unnecessary fields (insurance details, non-clinical identifiers)
- **Tiered Data Labeling and Access Control:** Implement role-based access control (RBAC) with granular permissions; encrypt patient data at rest (AES-256) and in transit (TLS 1.3)
- **Retention and De-identification:** Automatic purging of identifiable logs after 30 days; shifted to tokenized, de-identified audit trails for system performance monitoring
- **End-to-End Encryption:** Patient data encrypted before transmission to cloud services; decryption keys held by hospital IT, not third-party vendors

## **Organizational Controls**

- **Data Protection Officer (DPO):** Hospital-appointed DPO with authority to approve any new data fields or use cases before deployment
- **Quarterly Access Audits:** Automated log review to identify unauthorized access attempts; remediation within 48 hours
- **Default-Deny Data Sharing Policy:** Cross-project sharing of health data prohibited unless explicitly approved by DPO and hospital legal; stringent contractual language with third-party vendors

## **Societal/Compliance Controls**

- **Transparent Consent Mechanism:** Plain-language patient consent forms (written at 8th-grade reading level per NIH standards) clearly specifying data usage; patients can withdraw consent and request data export at any time during hospital admission
- **Breach Notification and Remediation:** Published breach notification protocol; patient notification within 72 hours of confirmed breach per HIPAA/PIPEDA; public transparency report on data incidents and remediation (quarterly publication)
- **Admission and Withdrawal:** Patients can withdraw from system at any point; patient opt-out recorded in EHR and respected across all hospital systems

## **Responsible Release Criteria & Metrics**

- **Unauthorized Access Incidents:** Target = 0 confirmed unauthorized access events per quarter
- **Overdue Data Retention:** Target = 100% compliance with 30-day purge policy; monthly audit confirmation

audit confirmation

- **Consent Withdrawal Fulfillment:** Target = 100% fulfillment within 2 business days of patient withdrawal request
- 

## *Mitigation 2: Diagnostic Error and Clinical Safety*

### **Technical Controls**

- **Uncertainty Quantification and High-Risk Escalation:** Model outputs include confidence scores and decision thresholds; clinical pathways designated as high-risk (chest pain, stroke, sepsis, obstetric emergencies) automatically escalate to mandatory physician review regardless of AI confidence score
- **Human Review Workflow:** System flags high-risk cases with confidence < 75% or cases meeting clinical criteria for serious pathology; these cases require physician sign-off before any clinical action
- **Pre-Deployment Simulation and Retrospective Validation:** System tested on de-identified historical cases with documented clinical outcomes; performance validated against ground-truth diagnoses confirmed by board-certified physicians

### **Organizational Controls**

- **Clinical Safety Officer:** Hospital-appointed Clinical Safety Officer (board-certified physician with patient safety expertise) approves all model updates before deployment
- **Monthly Postmortems:** Quality assurance review of diagnostic misses (false negatives in high-risk pathways); root cause analysis documented and fed into model retraining pipeline
- **Conservative A/B Testing Protocol:** New model versions tested only on low-risk pathways (e.g., routine musculoskeletal complaints) with documented safety thresholds before expansion

### **Societal/Compliance Controls**

- **Patient and Provider Communication:** Patient-facing UI displays disclaimer: *"This AI Assistant suggests possible diagnoses but is not a final diagnosis. A physician will review your case and confirm the diagnosis. Always seek emergency care for life-threatening symptoms."* Clinician UI displays confidence scores, supporting evidence, and alternative diagnoses; outputs presented with human-interpretable explanations rather than as black-box predictions

- **Informed Consent on Limitations:** Information sheet provided during hospital admission explaining the role of AI in triage and clinical decision-making; emphasis on physician authority in diagnosis

## Responsible Release Criteria & Metrics

- **100% Human Review Requirement:** All high-risk cases (chest pain, stroke, sepsis, obstetric) escalated to physician review and documented sign-off before clinical action
  - **Model-Attributable Serious Adverse Events (SAEs):** Target = 0 documented SAEs attributable to AI recommendation errors; quarterly trend analysis
  - **Clinician Override and Adjustment Rates:** Override rates monitored monthly; patterns indicating systematic model bias trigger retraining or threshold recalibration
- 

### *Mitigation 3: Fairness and Health Equity*

## Technical Controls

- **Stratified Accuracy Evaluation:** Model performance (accuracy, sensitivity, specificity, precision) evaluated separately by age group, gender, ethnicity, language, and socioeconomic indicators; performance gaps  $> 5\%$  trigger retraining or rule-based compensation
- **Disparity Trigger Mechanism:** System alerts if performance for any demographic group falls below target threshold; automatic escalation to fairness review board
- **Multilingual and Localized Interface:** System supports multiple languages (English, Spanish, Cantonese, Vietnamese, French); culturally contextualized symptom questioning (e.g., recognizing symptom descriptions unique to immigrant populations or cultural health beliefs)
- **Data Augmentation for Underrepresented Groups:** Active learning pipeline to systematically collect and label cases from underrepresented populations; quarterly data augmentation sprints

## Organizational Controls

- **Fairness Review Board:** Cross-functional board including physicians, nurses, patient advocates, representatives from minority communities, and data scientists; meets quarterly to review fairness metrics and recommend mitigations
- **Diverse Data Labeling and Annotation:** Human-labeling teams intentionally recruited from diverse backgrounds; bias-correction datasets explicitly built to address underrepresented groups
- **Targeted Data Collection:** Hospital quality improvement initiatives prioritize data collection from underrepresented patient populations (e.g., extending symptom intake for pediatric, geriatric, and non-English-speaking patients)

## Societal/Compliance Controls

- **Patient Feedback Mechanism:** In-app feedback channel allowing patients to report cases of misunderstood symptoms or non-applicable recommendations; feedback reviewed within 30 days
- **Public Fairness Reporting:** Annual fairness report published on hospital website documenting performance across demographic groups, identified gaps, and mitigation progress; transparent communication of remaining disparities

## Responsible Release Criteria & Metrics

- **Performance Gap Threshold:** Maximum performance gap (accuracy, sensitivity) across demographic groups = 5%; if exceeded, model update or rule adjustment triggered
  - **Feedback Closure Time:** 95% of patient feedback regarding system misunderstanding closed within 30 days; feedback patterns used to identify retraining opportunities
  - **Disparity Trend:** Year-over-year reduction in performance gaps across demographic groups; annual audits confirm progress
- 

## *Mitigation 4: Over-Reliance and Accountability*

## Technical Controls

- **Mandatory Display of Evidence and Alternatives:** System UI forces display of model confidence score, key decision factors (e.g., "presenting symptoms: fever + cough + hypoxia"), and alternative differential diagnoses; outputs presented with human-

interpretable explanations

- **Mandatory Human Confirmation Checkpoint:** Clinical decision writing to the electronic health record (EHR) requires explicit physician confirmation; system does not auto-populate clinical notes; physicians must manually enter their own diagnostic impression
- **Decision Audit Trail:** All AI recommendations and clinician responses (accept, modify, override) logged for retrospective auditing and feedback

## Organizational Controls

- **Clinician AI-Literacy Training:** Hospital-wide training program on AI-assisted decision-making, including psychology of automation bias, responsible use of decision support, and critical appraisal of AI outputs; annual refresher required for credentialing
- **Standard Operating Procedure (SOP) on AI Role:** Institutional policy explicitly states that physicians retain final diagnostic authority and legal responsibility for all clinical decisions; AI recommendations are advisory only; clear escalation pathway for cases where clinician disagrees with AI

## Societal/Compliance Controls

- **Patient Consent and Awareness:** Patients informed at admission that AI assists in their care; educational materials clarify that a clinician will make final diagnostic decisions
- **Public Institutional Statement:** Hospital publishes statement of AI role and limitations on website and in patient materials, addressing patient concerns about diagnostic authority and accountability

## Responsible Release Criteria & Metrics

- **100% Human Confirmation Rate:** All recommendations logged with corresponding clinician confirmation or override; 100% documentation rate maintained
- **Override/Modification Rate Stability:** Clinician override/modification rates monitored monthly; stable rates (10-20%) indicate appropriate use; elevated rates (>30%) trigger training review; low rates (<5%) may indicate insufficient clinician scrutiny
- **Complaint Resolution:** Institutional patient complaints regarding diagnostic misunderstandings tracked and declined year-over-year through education and interface improvements

---

## *Mitigation 5: Security and Attack Surface*

### **Technical Controls**

- **Threat Modeling and Adversarial Testing:** Formal threat model identifying attack vectors (model extraction, adversarial inputs, API abuse, ransomware); quarterly adversarial testing to identify vulnerabilities
- **Rate Limiting and Anomaly Detection:** API rate limiting (e.g., max 1000 requests/hour per user) prevents model extraction; anomaly detection monitors unusual query patterns (e.g., repeated similar queries with slight perturbations) and flags suspicious behavior
- **Key Rotation and Cryptographic Best Practices:** Encryption keys rotated quarterly; cryptographic algorithms updated to current standards per NIST recommendations
- **Incident Response Architecture:** Failover systems, rollback procedures, and feature toggles enable rapid response to security incidents; system can disable compromised features within minutes

### **Organizational Controls**

- **Incident Response Playbooks and Drills:** Documented incident response procedures for common attacks (model extraction, data breach, ransomware, DDoS); quarterly drills simulate security incidents
- **Third-Party Penetration Testing:** Annual independent security audits (OWASP Top 10 compliance); major security vulnerabilities fixed within 30 days of discovery
- **Disclosure and Remediation Protocol:** Major security incidents disclosed to hospital leadership and affected patients within 24 hours; public transparency report issued post-remediation

### **Societal/Compliance Controls**

- **Hospital Security Policy Alignment:** AI system architecture and security requirements aligned with hospital IT security policies, HIPAA security rule, and state/federal healthcare cybersecurity regulations
- **Patient Notification of Incidents:** If security incident occurs and patient data is compromised, patients notified within 72 hours per legal requirements; remediation plan communicated (e.g., credit monitoring, identity theft protection)

## Responsible Release Criteria & Metrics

- **Adversarial Test Coverage:** Target = 90% coverage of documented attack vectors through adversarial testing scenarios
  - **High-Severity Penetration Test Fixes:** Target = 100% of high-severity findings from pen tests remediated within 30 days
  - **Major Security Incidents:** Target = 0 confirmed major security incidents (data breach, unauthorized access, significant system compromise) per year
- 

## Unavoidable Biases

---

### *Bias 1: Representation Gaps in Training Data*

**Description:** Despite data augmentation and active learning efforts, clinical data are inherently imbalanced. Rare diseases, minority population presentations, and geographically remote patient populations are underrepresented in hospital datasets. Even with augmentation strategies, model confidence and accuracy lag on low-representation groups.

**Why Unavoidable:** Historical healthcare datasets reflect systemic inequities in healthcare access and research representation. Complete remediation requires decades of demographic-balanced data collection across all hospitals.

#### Developer Stance:

- Publish stratified performance metrics by demographic group in annual fairness reports; do not hide disparities
  - Implement default human-first pathway for low-representation groups: automatically escalate to senior physician review regardless of AI confidence
  - Maintain active learning roadmap with measurable targets for collecting additional diverse samples; partner with clinics serving underrepresented populations
  - Acknowledge limitations transparently to clinicians and patients
-

## *Bias 2: Distribution Shift and Data Drift*

**Description:** Patient symptom presentation, reporting styles, device quality, dialect, and vocabulary shift over time and across geographic regions. A model trained on 2024 data from urban academic hospitals may perform poorly in rural clinics or as disease epidemiology evolves (e.g., COVID-19 variants, emerging infections).

**Why Unavoidable:** Clinical environments are non-stationary; patient populations, disease prevalence, and presentation norms change continuously.

### **Developer Stance:**

- Continuous monitoring of input/output distribution: System tracks statistical properties of incoming symptom data; alert clinicians if distribution shifts significantly from training baseline
  - Drift detection triggers: When drift detected, recommend human review of high-risk cases and initiate limited retraining on recent data
  - Periodic retraining schedule: Formal retraining every 12 months or when drift alerts accumulate; retraining conducted with oversight from Clinical Safety Officer
- 

## *Bias 3: Inherent Value Trade-offs Between Sensitivity and Specificity*

**Description:** Clinical AI systems face an irreducible trade-off: maximizing sensitivity (catching all cases, reducing false negatives) increases false positives (unnecessary testing, cost, patient anxiety); prioritizing specificity (avoiding false alarms) increases false negatives (missed diagnoses, clinical harm).

**Why Unavoidable:** This is a fundamental property of binary classification. The "correct" balance is a value judgment reflecting institutional priorities, not an empirical fact discoverable through data alone.

### **Developer Stance:**

- Configurable thresholds governed by clinical leadership: Hospital clinical leadership and ethics board set decision thresholds for each clinical pathway, reflecting institutional risk tolerance
  - Default toward patient safety and higher recall: For life-threatening pathways (e.g., stroke, cardiac events), bias configuration toward high sensitivity (catch more cases even if false positives increase) to prioritize patient safety
  - Transparent documentation of trade-off choices: Publish the rationale for chosen thresholds and acknowledge limitations; allow settings adjustment by hospital if priorities shift
- 

## Conclusion

---

The Hospital AI Assistant presents significant potential for improving clinical efficiency, triage accuracy, and early detection of high-risk patients. However, it introduces substantive risks spanning data privacy, clinical safety, health equity, accountability, and cybersecurity—risks inherent to deploying AI in safety-critical healthcare environments.

The proposal implements a rigorous governance framework combining technical controls (data protection, explainability, access management), organizational oversight (Clinical Safety Officer, fairness review board, incident response), and societal accountability mechanisms (transparent reporting, patient consent, public fairness audits). These measures align with Microsoft Responsible AI Standards and reflect evidence from recent AI hospital deployments.

Residual risks remain unavoidable. Representation gaps in training data necessitate ongoing commitment to diverse data collection and conservative deployment in underrepresented groups. Distribution shift requires continuous monitoring and periodic model retraining to maintain performance as clinical contexts evolve. Trade-offs between sensitivity and specificity require ongoing governance and transparent documentation of institutional value choices.

Deployment should proceed as a controlled pilot in a single hospital department, with rigorous monitoring of clinical outcomes, adverse event tracking, and fairness metrics for 12 months before expansion. Annual impact assessments will inform refinements. Success depends on embedding human clinical judgment at the center of all high-consequence decisions and maintaining transparent communication with patients, clinicians, and the public regarding system capabilities and limitations.

---

## References

---

- Med-Tech World. (2025). The world's first AI hospital, developed in China, is transforming healthcare. Retrieved from <https://med-tech.world/news/chinas-ai-hospital-transforming-healthcare/>
- Microsoft. (2022a). *Microsoft Responsible AI Standard v2: General Requirements*. Retrieved from [http://www.microsoft.com/en-us/ai/principles-and-approach](https://www.microsoft.com/en-us/ai/principles-and-approach)
- Microsoft. (2022b). *Microsoft Responsible AI Impact Assessment Guide*. Retrieved from <https://www.microsoft.com/en-us/ai/principles-and-approach>
- Microsoft. (2022c). *Microsoft Responsible AI Impact Assessment Template*. Retrieved from <https://www.microsoft.com/en-us/ai/principles-and-approach>
- National Institute of Standards and Technology. (2023). AI Risk Management Framework (AI RMF 1.0). Retrieved from <https://www.nist.gov/itl/ai-risk-management-framework>
- Tsinghua University. (2025). Tsinghua AI Agent Hospital Inauguration and 2025 Tsinghua Medicine Townhall Meeting. Retrieved from <https://www.tsinghua.edu.cn/>