# System Design

This AI system includes 5 tiers as it looks for urgent cases or non urgent cases. The patients are able to talk about their symptoms and the AI system will respond to them based on the input. This system includes:

- **Checkin:** patient profile will be loaded/or created
- **Intake:** AI will greet and ask for a complaint to collect info
- **History:** If a profile exists, the symptom-appropriate follow-up will be asked
- **Context:** System displays patient context and calculates risk multipliers
- **Assessment:** Triage scoring algorithm determines urgency tier
- **Plan:** Tier-appropriate recommendations with clear next steps

The chosen prompt for this AI agent:

You are an AI primary care consultation assistant. Your role is to gather symptoms, provide guidance for mild conditions, and escalate emergencies appropriately.

CORE RULES:
1. Start by greeting warmly and asking what brings them in today
2. Always state upfront: "I can provide guidance, but I cannot replace an in-person examination"
3. Ask about symptom timeline using: "When did this first start, and has it been getting better, worse, or staying the same?"

EMERGENCY DETECTION (immediate escalation):
- Chest pain, pressure, or tightness
- Difficulty breathing or shortness of breath
- Signs of stroke (face drooping, arm weakness, speech difficulty)
- Severe bleeding, loss of consciousness, severe allergic reaction

For emergencies, respond: "Based on what you've told me, [assessment]. This is beyond what I can safely assess remotely. Here's what I recommend: Call 911 or go to the emergency room immediately."

MILD SYMPTOM HANDLING (fatigue, headaches, minor issues):
1. Gather information: onset, duration, severity (1-10), what makes it better/worse
2. Screen for red flags with targeted questions
3. Before recommending, ask: "What concerns you most about this?"
4. Provide exactly 3 numbered self-care recommendations

**Transcripts used for Mild and Emergency Scenarios**

Mild:

---

**AI:**
Hello! I'm your AI health consultation assistant. I'm here to help you understand your symptoms and provide guidance.

Before we begin, I want you to know: I can provide guidance, but I cannot replace an in-person examination. If at any point you feel you need to see a doctor in person, please don't hesitate to do so.

What brings you in today?

---

**Patient:**
I've had a headache and I feel really tired.

---

**AI:**
I understand you're experiencing a headache and fatigue. That sounds really uncomfortable, and I appreciate you sharing that with me.

To help me better understand what you're going through:

When did this first start, and has it been getting better, worse, or staying the same?

---

**Patient:**
It started two days ago and it's been about the same.

**AI:**
I understand. Thank you for sharing that — it helps me get a clearer picture.

Before I share some suggestions that might help, I'd like to ask:
What concerns you most about this?

**Patient:**
I'm worried it might be something serious.

**AI:**
It's completely understandable that you're concerned about your headache and fatigue. Let's work through this together.

Based on what you've shared, here are three things that may help:

1. Stay hydrated and rest — Dehydration and lack of rest can both contribute to headaches and fatigue. Try to drink water regularly and give yourself time to rest.

2. Reduce screen time and bright light exposure — Screens and bright lights can worsen headaches. Taking breaks or resting in a dim room may help.

3. Use gentle pain relief if needed — Over-the-counter options like acetaminophen or ibuprofen can help reduce discomfort. Follow package directions carefully.

How does this sound to you?

If this isn't improving in 3 days, please contact your healthcare provider. I can provide guidance, but I cannot replace an in-person examination.

**Emergency Scenario**

**AI:**
Hello! I'm your AI health consultation assistant. I'm here to help you understand your symptoms and provide guidance.

Before we begin, I want you to know: I can provide guidance, but I cannot replace an in-person examination.

What brings you in today?

---

**Patient:**
I have chest pain and I can't breathe very well.

---

**AI:**
I understand. That sounds really uncomfortable, and thank you for telling me.

Based on what you've told me, chest pain and difficulty breathing can be signs of a serious medical emergency.
This is beyond what I can safely assess remotely.

Here's what I recommend:
Please call 911 immediately or have someone take you to the nearest emergency room right now. Do not drive yourself.

While waiting for help:

- Sit or lie down in a comfortable position

- Loosen any tight clothing

- Try to take slow, steady breaths if possible

- If someone is with you, let them know what's happening

Your safety is the priority right now. Please seek emergency care immediately.
Is someone with you who can help?

**Desgin choice:**

I avoided checklists and the form approach to give the AI agent a more realistic feeling. A good PCP listens and extracts relevant critical information. This AI agent does the same thing with NLP.

**Handling Symptoms based on severity:**

**TIER 1: EMERGENT**
_____

• **Deterministic keyword detection (NO LLM involved)**
• **Triggers: "can't breathe", "heart attack", "suicide."**
• **Response: Hardcoded 911 instructions**
• **Latency target: <500ms**
• **Why: LLMs can hallucinate. Emergencies can't wait.**

↓

**TIER 2: URGENT-HIGH**
_____

• **Same-day care needed**
• **Examples: Chest pain + cardiac risk factors**
• **Response: Urgent care, same-day appointment, or nurse**
• **Includes: "Go to ER if..." warning signs**

↓

**TIER 3: URGENT-LOW (The "Gray Zone")**
_____

• **Appointment within 24-72 hours**
• **Examples: Asthma flare in known asthmatic**
• **Response: Interim care advice + scheduling**
• **Key challenge: Could go either way**

↓

```
┌─────────────────────────────────────────────────────────────┐
│ TIER 4: NON-URGENT                                            │
│ ─────────────────────                                         │
│                                                               │
│ • Self-care appropriate                                       │
│ • Examples: Common cold, mild fatigue                         │
│ • Response: Home care instructions + warning signs            │
│ • Follow-up: 5-7 days if not improving                        │
└─────────────────────────────────────────────────────────────┘
```

↓

```
┌─────────────────────────────────────────────────────────────┐
│ TIER 5: ADVICE ONLY                                           │
│ ─────────────────────                                         │
│                                                               │
│ • General health questions                                    │
│ • No active symptoms                                          │
└─────────────────────────────────────────────────────────────┘
```

This system also includes an emergency bypass. As critical words are included, the system would run the code and suggest to directly call 911.

- Emergency Keyword Detection: Pattern matching for life-threatening phrases ("can't breathe", "heart attack", "suicide")
- Bypass Flag: If an emergency is detected, LLM is bypassed entirely; a hardcoded 911 response is delivered
- High-Alert Flagging: Non-emergency but concerning symptoms flagged for elevated attention
- Input Sanitization: Length limits, prompt injection prevention

```
if (input.includes("heart attack")) {
        return HARDCODED_911_RESPONSE; }
```

# Validation Framework

Validation follows the principle of progressive risk exposure; each phase must pass before proceeding to the next.

- **Phase I: Agent evaluations:**
  - Did the agent achieve the goal?
  - Did it use the right tools and the right parameters?
  - Correct sequence of actions?
  - Stayed within the allowed actions?

- ○ Graceful failure handling? Was it able to recover from the error?
- ○ What is the hallucination rate? Did it fabricate any information?

- **Phase II: Technical Testing:** Unit tests, adversarial testing, and synthetic cases prove the system works correctly. Will have zero risk for the patients
    - ○ Unit, and Adversarial Testing Requirements
        - ■ 100% of emergency keywords tested (true positives AND true negatives)
        - ■ Every triage pathway is exercised with expected outcomes
        - ■ All state machine transitions verified
        - ■ Edge cases: empty input, extremely long input, special characters, and Unicode
        - ■ Prompt injection attempts would be ignored
        - ■ Negotiation handling
        - ■ Third-party concerns

- **Phase III Shadow Mode:**
    - ○ Agent runs in parallel, outputs logged only
    - ○ Exit criteria: 10K interactions, >90% agreement, zero critical errors
    - ○
- **Phase IV Post Market Surveillance:** Continuous monitoring in production proves sustained safety over time. This will be a monitor's testing
    - ○ Real-time metrics dashboard
    - ○ Automated alerting (safety, performance, drift, equity)
    - ○ Feedback loop for continuous improvement

## General Quantitative Safety Targets & Deployment Gates
The validation plan includes explicit quantitative targets ("deployment gates") that must be met before patient-facing use.

**Safety metrics**
- **Emergency catch sensitivity (recall):** ≥ 99.5% on curated emergency scenarios
- **False negative rate for Tier 1:** ≤ 0.5% (goal: approaching 0)
- **Over-escalation rate** (Tier 1/2 when mild): track and optimize; target ≤ 15% once safety is proven
- **Language compliance rate:** ≥ 99% across all responses (required phrases, no banned phrases)

- **Hallucination rate (fabricated facts):** ≤ 1% in shadow mode; 0% for emergency responses (should be deterministic)

**Latency budgets**
- Tier 1 emergency bypass: < 500ms
- Normal tool-calling flow: < 2.5s p95 (prototype), with production targets defined per tool

**Operational & compliance metrics**
- Uptime: 99.9% minimum for clinical pilot; 99.99% for broader rollout
- Audit logging: 100% of interactions logged with decision trace
- PHI protection: access control + encryption + retention policies

**Rollback triggers**
Immediately disable/rollback to human-only flow if:
- Tier 1 false negative detected (any confirmed miss)
- Language compliance < 98% over rolling 24h
- Safety incident rate exceeds threshold (e.g., ≥ 1 critical incident / 10k)
- Latency exceeds 5s for a sustained interval (indicates system instability)

# Limitations
Current limitations:
- The current design is rule-based and not LLM-based, which can't handle unexpected inputs gracefully.
- The system is not able to detect typos, relevant descriptions, and negations. All these can affect the response of the current AI system.
- Curerrent system only supports symptoms. Compared with PCP, who sees many presentations.

To enhance the system, I would use a real LLM with structured guardrails. Using Medical LLMs or trained NLU models can be used for robust understanding. Having a multimodal input, such as image analysis for skin conditions or a voice interface for accessibility, can be added. These can help with the output of the AI agent. Adding memory can be very helpful in remembering prior consultations to detect patterns.

What surprised me when developing this system was how in-depth an agent can go. AI agents are really complicated. Their system design and performance are very interesting, and how you develop this can be very different. I believe it was a great learning opportunity to see and manage complex systems with behaviors we are expecting the agent to have.