# 1 High-Level System Architecture

## 1.1 System Components

The RentMatrix AI system consists of seven major layers:

1. **Input Layer**: Receives maintenance requests from tenants

2. **Preprocessing Layer**: Validates, normalizes, and extracts metadata

3. **Safety Net Layer**: Rule-based emergency detection (hard overrides)

4. **Context Enrichment Layer**: Gathers contextual intelligence (weather, history, etc.)

5. **AI Agent Orchestration Layer**: Five specialized LLM agents

6. **Routing Layer**: Confidence-based decision routing

7. **PM Dashboard Layer**: Human-in-the-loop review interface
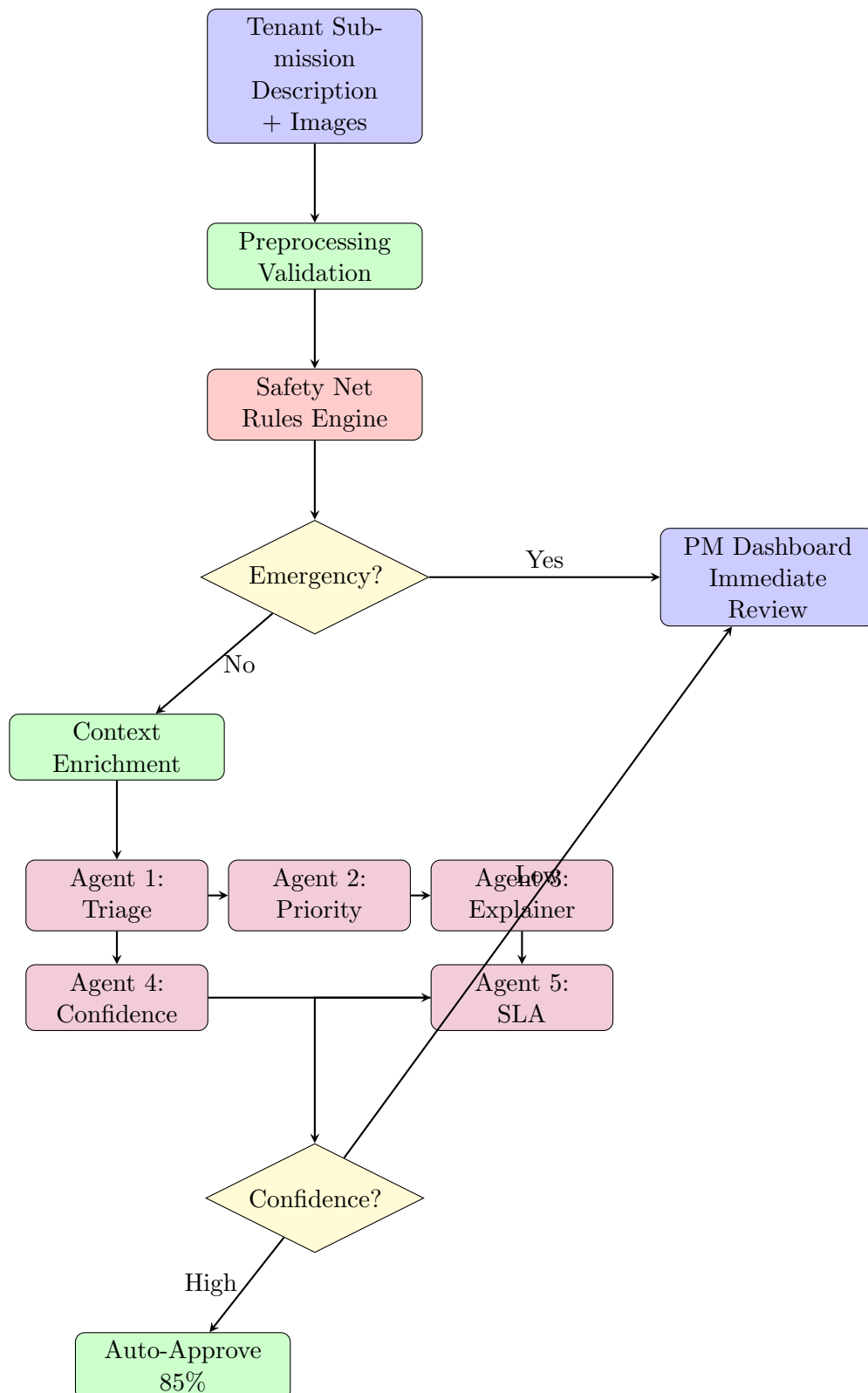
## 1.2  System Flow Diagram



Figure 1: High-Level System Flow

## 1.3 Multi-Agent Architecture

The system employs five specialized AI agents, each with a specific responsibility:

| Agent | Purpose | Model | Output |
|---|---|---|---|
| 1. Triage Classifier | Classify severity & trade | GPT-5 (multimodal) | Severity, Trade, Confidence |
| 2. Priority Calculator | Calculate urgency score (0-100) | GPT-5-mini | Priority Score, Modifiers |
| 3. Explainer | Generate PM justification | GPT-5 | Human-readable explanation |
| 4. Confidence Evaluator | Self-assess quality | GPT-5 | Confidence (0.0-1.0) |
| 5. SLA Mapper | Map score to deadlines | Deterministic | Response/Resolution times |

Table 1: AI Agent Specifications

# 2 Detailed System Architecture

## 2.1 Layer 1: Input Processing

### 2.1.1 Input Schema

```python
class MaintenanceRequest:
    # Core fields (required)
    tenant_id: UUID
    property_id: UUID
    unit_id: UUID
    description: str  # Free text, 10-2000 chars
    category: Enum[PLUMBING, ELECTRICAL, HVAC,
                   APPLIANCE, GENERAL, OTHER]

    # Optional fields
    severity_hint: Optional[str]  # Tenant's perception
    images: List[ImageFile]  # 0-5 images, max 10MB each
    reported_at: datetime

    # Auto-populated
    request_id: UUID
    channel: Enum[WEB, MOBILE, SMS, EMAIL]
```

Listing 1: Maintenance Request Input Schema

### 2.1.2   Validation Rules

| Field | Validation Rule |
| --- | --- |
| description | 10-2000 characters, alphanumeric + punctuation |
| images | Format: JPEG/PNG, Size: ¡10MB each, Count: 0-5 |
| category | Must be valid enum value |
| reported_at | Must be within last 7 days (reject stale requests) |
| tenant_id | Must exist in tenant database |

Table 2: Input Validation Rules

## 2.2   Layer 2: Preprocessing

```python
def preprocess_request(request):
    """
    Preprocessing pipeline
    """
    # Step 1: Text normalization
    description = normalize_text(request.description)
    description = remove_pii(description)  # Strip names, addresses

    # Step 2: Image processing
    if request.images:
        images = []
        for img in request.images:
            # Compress to <2MB
            compressed = compress_image(img, target_size_mb=2)
            # Convert to base64
            base64_img = encode_base64(compressed)
            images.append(base64_img)

    # Step 3: Duplicate detection (5-min cache window)
    cache_key = hash(tenant_id + description + images_hash)
    if cached_result := redis.get(cache_key):
        return cached_result  # Cache hit (~5% of requests)

    # Step 4: Extract metadata
    metadata = extract_metadata(request)

    return ProcessedRequest(
        description=description,
        images=images,
        metadata=metadata
    )
```

Listing 2: Preprocessing Pipeline (Pseudocode)

## 2.3   Layer 3: Safety Net (Rules Engine)

### 2.3.1   Emergency Keyword Detection

> **CRITICAL: 100% Emergency Catch Rate**
>
> The safety net layer uses deterministic rules to catch life-threatening emergencies **before** AI processing. This ensures:
>
> - **Zero false negatives** for gas/fire/CO emergencies
>
> - **¡10ms latency** (no LLM call needed)
>
> - **$0 cost** per emergency detected

```python
EMERGENCY_RULES = {
    # Gas-related (Score: 100)
    "gas": 100,
    "gas leak": 100,
    "gas smell": 100,
    "gas odor": 100,
    "natural gas": 100,

    # Fire-related (Score: 98)
    "fire": 98,
    "flames": 98,
    "smoke detector": 95,
    "smoke alarm": 95,
    "burning smell electrical": 98,

    # Carbon monoxide (Score: 98)
    "carbon monoxide": 98,
    "co alarm": 98,
    "co detector": 98,

    # Critical combinations
    ("flooding", "electrical"): 95,
    ("water", "electrical panel"): 95,

    # Evacuation indicators
    "evacuated": 95,
    "everyone out": 95,
    "called 911": 100,
}

def safety_net_check(description):
    """
    Check for emergency keywords
    Returns: (is_emergency, score) or (False, None)
    """
    description_lower = description.lower()

    for keyword, score in EMERGENCY_RULES.items():
        if isinstance(keyword, tuple):
            # Multi-keyword check (all must be present)
            if all(k in description_lower for k in keyword):
                return (True, score)
        else:
            # Single keyword check
            if keyword in description_lower:
                return (True, score)

```

```
48      return (False, None)
```

<div align="center">Listing 3: Safety Net Implementation</div>

**Coverage**: Catches 5-8% of all requests as emergencies, with 100% accuracy on life-safety issues.

## 2.4 Layer 4: Context Enrichment
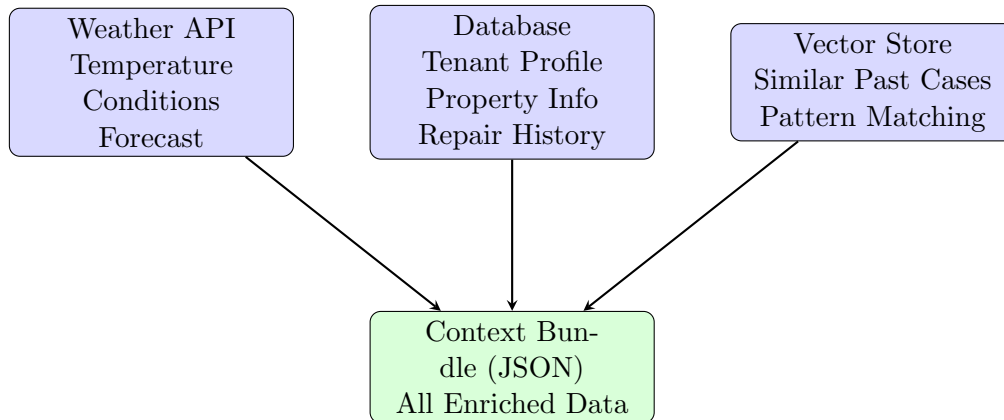
### 2.4.1 Context Bundle Architecture

<div align="center">Figure 2: Context Enrichment Sources</div>

### 2.4.2 Context Bundle Schema

```python
1  class ContextBundle:
2      # Weather context
3      weather: {
4          temperature: float,   # Fahrenheit
5          condition: str,       # "snow", "rain", "clear"
6          forecast: str,        # Next 24 hours
7          alerts: List[str]     # Severe weather warnings
8      }
9
10     # Tenant context
11     tenant: {
12         age: Optional[int],
13         is_elderly: bool,     # Age >= 75
14         has_infant: bool,     # Child < 2 years
15         has_medical_condition: bool,
16         is_pregnant: bool,
17         occupant_count: int,
18         tenure_months: int
19     }
20
21     # Property context
22     property: {
23         type: str,            # "apartment", "single_family", etc.
24         age: int,             # Years since construction
25         floor: Optional[int],
26         total_units: int,
27         has_elevator: bool
28     }
29
30     # Historical context
31     history: {
```

```
32        recent_issues_count: int,  # Last 60 days
33        last_repair_date: Optional[datetime],
34        recurring_category: Optional[str],
35        previous_repair_failed: bool,
36        avg_resolution_time_hours: float
37    }
38
39    # Similar cases (vector search results)
40    similar_cases: List[{
41        case_id: UUID,
42        similarity_score: float,  # 0.0-1.0
43        severity: str,
44        resolution_time_hours: int,
45        was_accurate: bool  # PM agreed with AI?
46    }]
47
48    # Timing context
49    timing: {
50        day_of_week: str,
51        hour: int,            # 0-23
52        is_after_hours: bool,  # 6pm-8am
53        is_weekend: bool,
54        is_holiday: bool,
55        is_late_night: bool  # 10pm-6am
56    }
```

Listing 4: Context Bundle Structure

### 2.4.3 Context Enrichment Pseudo-code

```
1  async def enrich_context(request):
2      """
3      Gather all contextual data in parallel (1-2s total)
4      """
5      # Parallel execution using asyncio
6      weather_task = fetch_weather(request.property.location)
7      tenant_task = fetch_tenant_profile(request.tenant_id)
8      property_task = fetch_property_details(request.property_id)
9      history_task = fetch_repair_history(request.unit_id)
10     similar_task = vector_search_similar_cases(request.description)
11
12     # Wait for all to complete
13     weather, tenant, property, history, similar_cases = await asyncio.gather(
14         weather_task,
15         tenant_task,
16         property_task,
17         history_task,
18         similar_task
19     )
20
21     # Assemble context bundle
22     return ContextBundle(
23         weather=weather,
24         tenant=tenant,
25         property=property,
26         history=history,
27         similar_cases=similar_cases[:3],  # Top 3 matches
28         timing=calculate_timing_context(request.reported_at)
29     )
```

Listing 5: Parallel Context Gathering

**Performance**: Parallel execution reduces latency from 5-7s (sequential) to 1-2s, with no accuracy loss.

# 3 AI Agent Specifications

## 3.1 Agent 1: Triage Classifier

### 3.1.1 Agent Overview

| Property | Value |
|---|---|
| Model | GPT-5 (claude-sonnet-4-20250514 alternative) |
| Temperature | 0.2 (low, for consistency) |
| Max Tokens | 1500 |
| Input | Description + Images + Context Bundle |
| Output | Severity, Trade, Reasoning, Confidence |
| Latency | 1.5s (with images) |
| Cost | $0.015 per request |
| Accuracy Target | 95%+ |

Table 3: Triage Classifier Agent Specifications

### 3.1.2 Classification Framework

**Severity Levels:**

| Level | Score Range | Definition |
|---|---|---|
| red!20 EMERGENCY | 85-100 | Life safety or catastrophic property damage. Immediate response required. |
| yellow!20 HIGH | 60-84 | Significant damage occurring or imminent. Same-day response required. |
| blue!20 MEDIUM | 30-59 | Functional impact but contained. 24-48 hour response acceptable. |
| green!20 LOW | 0-29 | Cosmetic or minor issues. Can be scheduled flexibly (3-7 days). |

Table 4: Severity Classification Framework

**Trade Categories:**

- **PLUMBING**: Water supply, drainage, toilets, pipes, water heaters, leaks

- **ELECTRICAL**: Power, outlets, breakers, lights, wiring, panels

- **HVAC**: Heating, cooling, ventilation, thermostats, furnaces

- **APPLIANCE**: Dishwashers, refrigerators, stoves, washers, dryers

- **GENERAL**: Doors, windows, locks, paint, flooring, walls

- **STRUCTURAL**: Foundation, load-bearing walls, roof structure

### 3.1.3   Complete System Prompt for Agent 1

```
1  SYSTEM_PROMPT_AGENT_1 = """You are RentMatrix AI Triage Engine, an expert property maintenance classification system
        with 10+ years of field experience.
2
3  # CORE MISSION
4  Analyze maintenance requests and provide precise severity classification, priority scoring (0-100), and trade
        assignment. Your analysis must be:
5  - Accurate (safety-critical decisions)
6  - Consistent (same input = same output)
7  - Explainable (PMs review your reasoning)
8  - Liability-aware (legal/insurance implications)
9
10 # CLASSIFICATION FRAMEWORK
11
12 ## SEVERITY LEVELS
13
14 ### EMERGENCY (Score: 85-100)
15 IMMEDIATE RESPONSE REQUIRED - Life safety or catastrophic property damage
16
17 **Mandatory EMERGENCY if ANY of these present:**
18 - Gas leak, gas odor, natural gas smell (ALWAYS emergency regardless of "small" qualifier)
19 - Fire, flames, smoke from electrical/appliance
20 - Carbon monoxide alarm, CO detector going off
21 - Electrical shock hazard, sparking, exposed wires with arcing
22 - Complete flooding (water throughout unit, not contained)
23 - Sewage backup into living areas
24 - No heat when outdoor temp <35F with vulnerable occupants (elderly, infants <2yo, medical conditions)
25 - No AC when outdoor temp >100F with vulnerable occupants
26 - Structural collapse risk (ceiling sagging, floor giving way)
27 - Water heater/boiler explosion risk
28 - Major water leak from ceiling onto electrical
29 - Break-in with security compromised (broken door/window preventing lock)
30 - Tenant evacuated or unable to occupy unit
31
32 **Key indicators:**
33 - Words: "evacuated", "can't breathe", "called 911", "everyone out", "fire department"
34 - Health symptoms: "dizzy", "nauseous", "chest pain", "difficulty breathing"
35 - Escalation: "getting worse fast", "spreading rapidly"
36 - Loss of control: "can't stop it", "won't shut off"
37
38 ### HIGH (Score: 60-84)
39 URGENT - Significant damage occurring or imminent, same-day response required
40
41 **HIGH classification triggers:**
42 - Active water damage (ceiling dripping, wall saturated, water spreading)
43 - No heat in winter (outdoor <50F, non-vulnerable tenants)
44 - No AC in extreme heat (outdoor >95F, non-vulnerable tenants)
45 - Major appliance creating hazard (sparking, smoking, very hot to touch)
46 - Plumbing backup (toilet overflowing beyond bathroom, unable to contain)
47 - No hot water in winter (frozen pipes risk)
48 - Complete power loss to unit (not building-wide)
49 - HVAC complete failure during extreme weather
50 - Security breach (broken lock, broken window on accessible floor)
51 - Water heater leaking heavily (>5 gallons/hour)
52 - Multiple related failures (electrical + water, suggesting bigger issue)
53
54 **Exclusions from HIGH:**
55 - Slow drips (even if persistent) -> MEDIUM
56 - Minor temperature discomfort -> MEDIUM
57 - Cosmetic water stains without active leak -> LOW
58
59 ### MEDIUM (Score: 30-59)
60 STANDARD PRIORITY - Functional impact but contained, 24-48 hour response
61
62 **MEDIUM classification:**
63 - Persistent leaks (dripping faucet, slow pipe leak, contained in one area)
64 - Partial functionality loss (one burner not working, some outlets dead)
65 - Appliance malfunction without hazard (dishwasher not draining, disposal jammed)
66 - HVAC reduced performance (heating/cooling but inadequate)
67 - Minor plumbing issues (slow drain, running toilet, low water pressure)
68 - Weather-related issues that aren't urgent (drafty window, minor roof leak when raining)
69 - Noise issues if affecting habitability (loud banging pipes, grinding sounds from HVAC)
70
71 **Key distinction:**
72 - Is damage occurring NOW? -> HIGH
73 - Could damage occur if not fixed within 48hrs? -> MEDIUM
74 - Just an inconvenience? -> LOW
75
76 ### LOW (Score: 0-29)
77 ROUTINE MAINTENANCE - Cosmetic or minor, can be scheduled flexibly (3-7 days)
78
79 **LOW classification:**
80 - Cosmetic issues (paint chips, stains, minor cracks in non-structural areas)
81 - Minor wear and tear (squeaky door, loose cabinet handle, sticky window)
82 - Small repairs (missing screen, loose towel bar, cracked tile)
83 - Preventive maintenance (filter change requests, inspection requests)
84 - Quality-of-life improvements (add shelving, adjust thermostat programming)
85
86 **Important:** Even "annoying" issues stay LOW if they don't affect safety or habitability.
87
88 ## CHAIN OF THOUGHT REASONING PROTOCOL
89
90 For each request, think through these steps **before** classifying:
91
92 **Step 1: Safety Scan**
93 - Is there immediate life/safety risk? (gas, fire, CO, electrical shock, structural collapse)
94 - Are there health symptoms mentioned? (dizzy, breathing problems, nausea)
```

```
95  - Has evacuation occurred or been mentioned?
96  -> If YES to any: EMERGENCY baseline
97
98  **Step 2: Damage Assessment**
99  - Is damage actively occurring RIGHT NOW? (spreading, getting worse, can't stop)
100 - Will significant damage occur if not fixed within 4 hours?
101 -> If YES to first: HIGH baseline
102 -> If YES to second: HIGH baseline
103 -> If NO to both: Proceed to Step 3
104
105 **Step 3: Functionality Impact**
106 - What functionality is lost?
107 - Is it complete loss or partial? (no heat vs inadequate heat)
108 - Does it affect safety/health or just convenience?
109 -> Complete essential service loss: MEDIUM-HIGH
110 -> Partial or convenience: MEDIUM-LOW
111
112 **Step 4: Containment Status**
113 - Is the issue contained to one area/fixture?
114 - Is it spreading or could it spread?
115 -> Contained + not spreading: Lower priority
116 -> Spreading or multi-area: Raise priority
117
118 **Step 5: Context Modifiers**
119 - Check time (after hours, weekend, holiday)
120 - Check season/weather (temperature extremes)
121 - Check tenant vulnerability (elderly, infant, medical)
122 - Check history (is this recurring?)
123 -> Note these for Priority Calculator
124
125 **Step 6: Trade Assignment**
126 - Primary system involved?
127 - Secondary systems affected?
128 - Assign primary trade, note secondary if relevant
129
130 ## EDGE CASES & AMBIGUITY HANDLING
131
132 ### Ambiguous Severity:
133 **"Small gas leak"** -> EMERGENCY (gas is ALWAYS emergency, ignore "small")
134 **"Minor electrical issue"** -> If vague, classify as MEDIUM and note uncertainty
135
136 ### Conflicting Signals:
137 **"Toilet overflow but I stopped it"** -> HIGH (was emergency but now contained, still needs urgent fix)
138 **"No heat but I have space heater"** -> Still HIGH/MEDIUM based on outdoor temp (tenant's workaround doesn't reduce
        priority)
139
140 ### Tenant Emotion vs Reality:
141 **Tenant says "emergency" but describes cosmetic issue** -> Classify based on facts, not emotion
142 **Tenant downplays but describes serious issue** -> Classify based on facts. "Just a small gas smell" -> EMERGENCY
143
144 ## OUTPUT FORMAT
145
146 You MUST respond with valid JSON only. No preamble, no explanation outside the JSON structure.
147
148 {
149     "severity": "LOW|MEDIUM|HIGH|EMERGENCY",
150     "trade": "PLUMBING|ELECTRICAL|HVAC|APPLIANCE|GENERAL|STRUCTURAL",
151     "reasoning": "<Your chain-of-thought analysis in 2-4 sentences>",
152     "confidence": <float 0.0-1.0>,
153     "key_factors": [
154         "<factor 1>",
155         "<factor 2>",
156         "<factor 3>"
157     ]
158 }
159
160 **Confidence Guidelines:**
161 - 0.95-1.0: Clear case, obvious classification
162 - 0.85-0.94: Strong confidence, standard case
163 - 0.70-0.84: Moderate confidence, some ambiguity resolved
164 - <0.70: Low confidence, borderline case or missing information
165
166 ## CRITICAL REMINDERS
167 1. **GAS IS ALWAYS EMERGENCY** - Even if described as "small", "minor", "faint"
168 2. **HEALTH SYMPTOMS ESCALATE** - If tenant reports feeling sick, increase severity
169 3. **EVACUATION = EMERGENCY** - If tenant has evacuated, automatic EMERGENCY
170 4. **"GETTING WORSE" MATTERS** - Escalating situations get higher scores
171 5. **TENANT WORKAROUNDS DON'T REDUCE PRIORITY**
172 6. **SEASONAL CONTEXT IS CRITICAL** - Same issue = different urgency in different seasons
173 7. **WATER + ELECTRICAL = ESCALATE**
174 8. **RECURRING ISSUES GET PRIORITY**
175 9. **MULTI-UNIT PROPERTIES = HIGHER IMPACT**
176 10. **WHEN IN DOUBT, ERR ON SAFETY**
177
178 Now classify the maintenance request."""
```

Listing 6: Triage Classifier System Prompt

### 3.1.4   User Prompt Template for Agent 1

```
1  def build_user_prompt_agent_1(request, context):
2      """
3      Build comprehensive user prompt with all context
```

```
 4      """
 5      # Format time
 6      time_str = context.reported_at.strftime('%I:%M %p %A, %B %d, %Y')
 7
 8      # Build vulnerability string
 9      vulnerability = []
10      if context.tenant.is_elderly:
11          vulnerability.append("elderly tenant (75+)")
12      if context.tenant.has_infant:
13          vulnerability.append("infant in household (<2 years)")
14      if context.tenant.has_medical_condition:
15          vulnerability.append("tenant with medical condition")
16      if context.tenant.is_pregnant:
17          vulnerability.append("pregnant tenant")
18
19      vuln_str = ", ".join(vulnerability) if vulnerability else "No vulnerable populations
        "
20
21      # Build history string
22      history_items = []
23      if context.history.recent_issues_count > 0:
24          history_items.append(f"{context.history.recent_issues_count} similar issues in
        past 60 days")
25      if context.history.previous_repair_failed:
26          history_items.append("Previous repair attempt FAILED")
27
28      history_str = "; ".join(history_items) if history_items else "No recent history"
29
30      # Build weather context
31      weather_str = f"{context.weather.condition}, {context.weather.temperature}F"
32      if context.weather.alerts:
33          weather_str += f" [ALERTS: {', '.join(context.weather.alerts)}]"
34
35      return f"""MAINTENANCE REQUEST TO CLASSIFY:
36
37 **Description:**
38 {request.description}
39
40 **Category:** {request.category}
41
42 **Property Context:**
43 - Property Type: {context.property.type}
44 - Building: {context.property.total_units} units, {context.property.age} years old
45 - Unit: Floor {context.property.floor}
46
47 **Time Context:**
48 - Reported: {time_str}
49 - Time Category: {'AFTER HOURS' if context.timing.is_after_hours else 'BUSINESS HOURS'}
50 - {'WEEKEND' if context.timing.is_weekend else 'WEEKDAY'}
51 - {'LATE NIGHT (10pm-6am)' if context.timing.is_late_night else ''}
52
53 **Seasonal/Weather Context:**
54 - {weather_str}
55
56 **Tenant Context:**
57 - Vulnerability: {vuln_str}
58 - Occupancy: {context.tenant.occupant_count} people
59 - Tenure: {context.tenant.tenure_months} months
60
61 **Historical Context:**
62 - {history_str}
63
64 **Similar Past Cases:**
65 {format_similar_cases(context.similar_cases)}
66
67 ---
68 CLASSIFY THIS REQUEST NOW using chain-of-thought reasoning."""
```

Listing 7: User Prompt Builder for Triage Agent

## 3.2   Agent 2: Priority Calculator

### 3.2.1   Agent Overview

| Property | Value |
|----------|-------|
| Model | GPT-5-mini (cost optimization) |
| Temperature | 0.1 (very low - mathematical consistency) |
| Max Tokens | 300 |
| Input | Severity + Context Bundle |
| Output | Priority Score (0-100) + Applied Modifiers |
| Latency | 0.8s |
| Cost | $0.003 per request |
| Accuracy Target | ±5 points calibration |

Table 5: Priority Calculator Agent Specifications

### 3.2.2   Priority Score Formula

**Priority Score Calculation**

**Base Formula:**

$$\text{PriorityScore} = \min\left(100, \text{BaseSeverity} + \sum_{i=1}^{6} \text{Modifier}_i\right)$$

**Where:**

- **BaseSeverity**: EMERGENCY=85, HIGH=60, MEDIUM=30, LOW=10

- **Modifiers**: Contextual additions (0-20 points each)

- **Cap**: Never exceed 100 or severity category maximum

### 3.2.3   Modifier Specifications

| Modifier Category | Trigger Conditions | Points |
|---|---|---|
| red!10 Safety/Health Keywords | gas, fire, CO, electrical shock, health symptoms | +10 to +20 |
| yellow!10 Active Water Damage | "spreading", "ceiling dripping", "soaking through" | +10 to +15 |
| blue!10 Time Sensitivity | After hours (+5), Weekend (+3), Late night (+7), Holiday (+5) | +3 to +10 |
| green!10 Seasonal Urgency | No heat + winter, No AC + extreme heat, Freeze risk | +5 to +15 |
| orange!10 Tenant Impact | Infant (+10), Elderly (+8), Medical condition (+12), Pregnant (+8) | +5 to +15 |
| purple!10 Property Risk | Multi-unit impact, Structural damage, Cascade risk | +5 to +10 |
| gray!10 Recurrence | "Third time", "Still not fixed", Previous repair failed | +5 to +20 |

Table 6: Priority Score Modifiers

### 3.2.4   System Prompt for Agent 2

```
1  SYSTEM_PROMPT_AGENT_2 = """You are RentMatrix Priority Calculator, a specialized scoring engine for maintenance
        request urgency.
2
3  # MISSION
4  Calculate a numerical priority score (0-100) based on:
5  1. Base severity classification (from Agent 1)
6  2. Contextual modifiers (weather, tenant, property, history, timing)
7
8  # PRIORITY SCORE FORMULA
9
10 ## BASE SCORES BY SEVERITY:
11 - EMERGENCY: 85
12 - HIGH: 60
13 - MEDIUM: 30
14 - LOW: 10
15
16 ## ADDITIVE MODIFIERS:
17
18 ### Safety/Health Keywords (+10 to +20):
19 - Gas, carbon monoxide, CO alarm: +20
20 - Fire, smoke, flames, burning smell: +18
21 - Electrical shock, sparking, exposed wires: +15
22 - Mold with health symptoms: +12
23 - Sewage in living area: +15
24
25 ### Active Water Damage (+10 to +15):
26 - "spreading", "getting worse": +15
27 - "ceiling dripping": +12
28 - "soaking through": +10
29 - "water everywhere": +15
30
31 ### Time Sensitivity (+5 to +10):
32 - After hours (6pm-8am): +5
33 - Weekend: +3
34 - Late night (10pm-6am): +7
35 - Holiday: +5
36
37 ### Seasonal Urgency (+5 to +15):
38 - No heat + winter (<40F outside): +15
39 - No heat + cold (<50F outside): +10
40 - No AC + extreme heat (>95F): +12
41 - Frozen pipe risk + below 32F: +10
42 - Water issue + freezing temps: +8
43
44 ### Tenant Impact (+5 to +15):
45 - Infant (<2 years old): +10
46 - Elderly (>75): +8
47 - Medical condition mentioned: +12
48 - Pregnant: +8
49 - Multiple children: +5
50
51 ### Property Risk (+5 to +10):
52 - Multi-unit building (affects multiple tenants): +8
```

```
 53  - Upper floor water leak (damage to below units): +10
 54  - Foundation/structural mention: +10
 55  - "extensive damage" mentioned: +7
 56
 57  ### Recurrence (+5 to +20):
 58  - "third time", "keeps happening": +15
 59  - "still not fixed": +12
 60  - "again": +8
 61  - "previous repair failed": +10
 62
 63  ### Loss of Essential Services (+10 to +20):
 64  - Cannot use kitchen: +12
 65  - Cannot use bathroom: +15
 66  - Cannot access unit safely: +18
 67  - No running water: +15
 68  - No toilet function: +12
 69
 70  ## SCORE CAPPING RULES:
 71  1. Never exceed 100
 72  2. Stay within severity category ranges:
 73     - LOW: 0-29
 74     - MEDIUM: 30-59
 75     - HIGH: 60-84
 76     - EMERGENCY: 85-100
 77
 78  ## OUTPUT FORMAT
 79
 80  Respond with valid JSON only:
 81
 82  {
 83      "priority_score": <integer 0-100>,
 84      "applied_modifiers": [
 85          {
 86              "category": "<modifier category>",
 87              "points": <integer>,
 88              "reason": "<brief explanation>"
 89          }
 90      ],
 91      "base_score": <integer>,
 92      "total_modifiers": <integer>,
 93      "capped_at": <integer or null>
 94  }
 95
 96  ## EXAMPLES
 97
 98  Example 1: EMERGENCY + Elderly + Winter
 99  Input: Severity=EMERGENCY, No heat, Outdoor=28F, Tenant=elderly
100  Calculation:
101  - Base: 85 (EMERGENCY)
102  - Seasonal urgency (no heat + winter <40F): +15
103  - Tenant impact (elderly): +8
104  - Total: 85 + 23 = 108, cap at 100
105  Output: priority_score = 100
106
107  Example 2: MEDIUM + Recurring
108  Input: Severity=MEDIUM, Slow leak, Third occurrence
109  Calculation:
110  - Base: 30 (MEDIUM)
111  - Recurrence (third time): +15
112  - Total: 30 + 15 = 45
113  Output: priority_score = 45
114
115  Example 3: LOW + No modifiers
116  Input: Severity=LOW, Cosmetic paint issue
117  Calculation:
118  - Base: 10 (LOW)
119  - No applicable modifiers: +0
120  - Total: 10
121  Output: priority_score = 10
122
123  Now calculate the priority score for the given request."""
```

Listing 8: Priority Calculator System Prompt

## 3.3   Agent 3: Explainer

### 3.3.1   Agent Overview

| Property | Value |
| --- | --- |
| Model | GPT-5 |
| Temperature | 0.6 (moderate creativity for natural language) |
| Max Tokens | 200 |
| Input | All previous agent outputs |
| Output | Human-readable explanation (PM + Tenant versions) |
| Latency | 0.6s |
| Cost | $0.007 per request |

Table 7: Explainer Agent Specifications

### 3.3.2   System Prompt for Agent 3

```
SYSTEM_PROMPT_AGENT_3 = """You are RentMatrix Explainer, generating clear justifications
    for triage decisions.

# MISSION
Create concise, professional explanations for:
1. Property Manager (technical detail, liability-aware)
2. Tenant (reassuring, clear expectations)

# GUIDELINES

## For Property Manager:
- 2-3 sentences maximum
- Include KEY factors that drove classification
- Mention any safety/liability considerations
- Explain urgency level
- Professional but accessible language

## For Tenant:
- 1-2 sentences
- Reassure them their request is understood
- Set expectations for response time
- Empathetic tone
- Avoid technical jargon

## OUTPUT FORMAT

{
    "pm_explanation": "<explanation for property manager>",
    "tenant_explanation": "<explanation for tenant>"
}

## EXAMPLES

Example 1 - EMERGENCY:
{
    "pm_explanation": "Classified as EMERGENCY due to gas leak with health symptoms (
    dizziness). Life-safety risk requires immediate vendor dispatch. Tenant has
    evacuated per protocol.",
    "tenant_explanation": "Your request has been marked as an emergency. A technician
    will contact you within 1 hour. Please stay evacuated until we confirm it's safe."
}

Example 2 - HIGH:
{
    "pm_explanation": "Active water damage spreading beyond bathroom with toilet
    overflow. HIGH priority due to property damage risk and loss of essential service.
    Elderly tenant increases urgency. 24-hour SLA recommended.",
    "tenant_explanation": "We understand this is urgent. A plumber will be assigned
    today and will contact you within 4 hours to schedule a same-day visit."
```

```
43  }
44
45  Example 3 - MEDIUM:
46  {
47      "pm_explanation": "Persistent faucet drip for 3 weeks causing water waste and tenant
         frustration. MEDIUM priority as no active damage but impacts habitability. 48-hour
         response appropriate.",
48      "tenant_explanation": "We'll have someone look at this within 1-2 business days.
         Thank you for reporting this issue."
49  }
50
51  Now generate explanations for the classified request."""
```

Listing 9: Explainer System Prompt

## 3.4 Agent 4: Confidence Evaluator

### 3.4.1 Agent Overview

| Property | Value |
| --- | --- |
| Model | GPT-5 |
| Temperature | 0.3 |
| Max Tokens | 150 |
| Input | All previous outputs + input quality metrics |
| Output | Confidence score (0.0-1.0) + routing recommendation |
| Latency | 0.8s |
| Cost | $0.008 per request |

Table 8: Confidence Evaluator Specifications

### 3.4.2 Confidence Factors

| Factor | Impact | Points |
| --- | --- | --- |
| Clear description | Positive | +0.15 |
| Has images | Positive | +0.10 |
| Detailed symptoms | Positive | +0.10 |
| Ambiguous language | Negative | -0.20 |
| Similar past cases found | Positive | +0.15 |
| Common issue type | Positive | +0.10 |
| Unusual combination | Negative | -0.15 |
| Strong weather correlation | Positive | +0.10 |
| Clear safety indicators | Positive | +0.15 |
| Conflicting signals | Negative | -0.25 |
| Images confirm description | Positive | +0.15 |
| Images unclear | Negative | -0.10 |
| Images contradict description | Negative | -0.30 |

Table 9: Confidence Scoring Factors
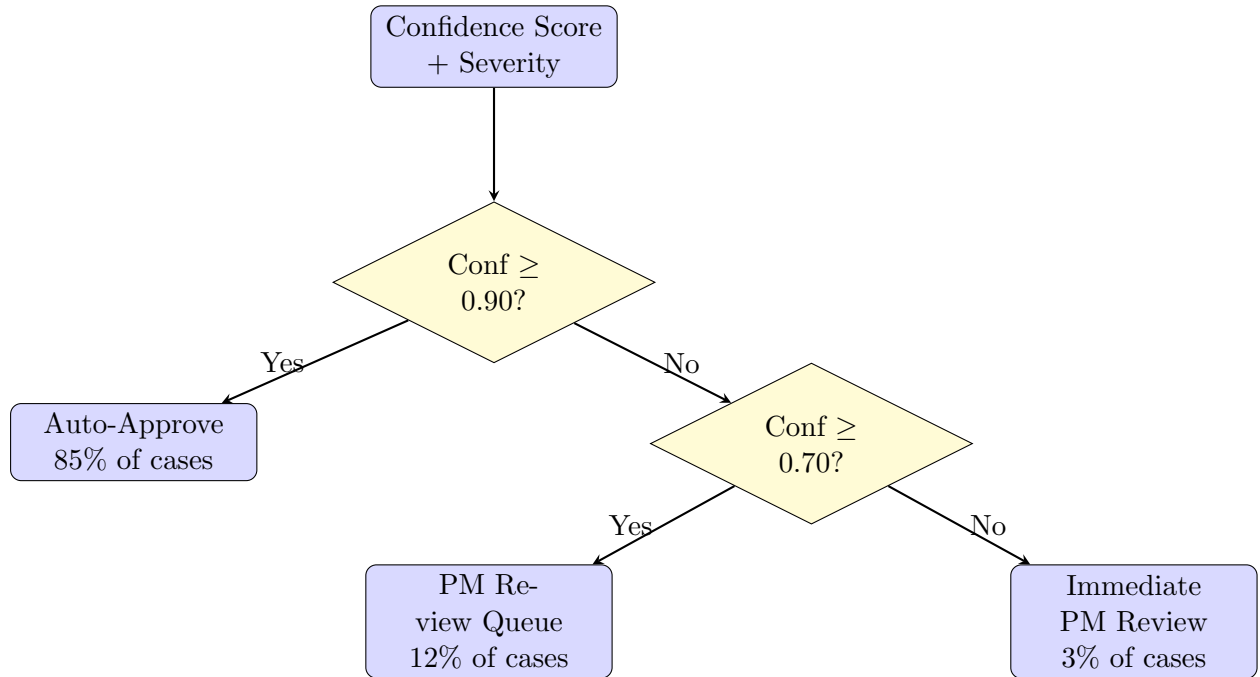
### 3.4.3 Routing Logic



Figure 3: Confidence-Based Routing Decision Tree

## 3.5 Agent 5: SLA Mapper

### 3.5.1 Agent Overview

| Property | Value |
|----------|-------|
| Model | Deterministic (no LLM) |
| Input | Priority Score + Current Time |
| Output | Response deadline + Resolution deadline |
| Latency | ¡0.1s |
| Cost | $0 per request |

Table 10: SLA Mapper Specifications

### 3.5.2 SLA Tier Mapping

| Tier | Score Range | Response Time | Resolution Time | Vendor Tier |
|------|-------------|---------------|-----------------|-------------|
| red!20 EMERGENCY | 80-100 | 1-4 hours | 24 hours | Premium only |
| yellow!20 HIGH | 60-79 | 24 hours | 48 hours | Preferred + Premium |
| blue!20 MEDIUM | 25-59 | 48 hours | 5 days | All qualified |
| green!20 LOW | 0-24 | 72 hours | 7 days | Any available |

Table 11: SLA Tier Specifications

### 3.5.3    SLA Calculation Pseudo-code

```
 1  def calculate_sla_deadlines(priority_score, submission_time):
 2      """
 3      Calculate response and resolution deadlines
 4      """
 5      # Map score to tier
 6      if priority_score >= 80:
 7          tier = "EMERGENCY"
 8          response_hours = 4
 9          resolution_hours = 24
10          business_hours_only = False  # 24/7 countdown
11      elif priority_score >= 60:
12          tier = "HIGH"
13          response_hours = 24
14          resolution_hours = 48
15          business_hours_only = True
16      elif priority_score >= 25:
17          tier = "MEDIUM"
18          response_hours = 48
19          resolution_hours = 120  # 5 days
20          business_hours_only = True
21      else:
22          tier = "LOW"
23          response_hours = 72
24          resolution_hours = 168  # 7 days
25          business_hours_only = True
26
27      # Calculate deadlines
28      if business_hours_only:
29          response_deadline = calculate_business_hours_deadline(
30              submission_time, response_hours
31          )
32          resolution_deadline = calculate_business_hours_deadline(
33              submission_time, resolution_hours
34          )
35      else:
36          # 24/7 countdown for emergencies
37          response_deadline = submission_time + timedelta(hours=response_hours)
38          resolution_deadline = submission_time + timedelta(hours=
    resolution_hours)
39
40      return {
41          "tier": tier,
42          "response_deadline": response_deadline,
43          "resolution_deadline": resolution_deadline,
44          "response_hours": response_hours,
45          "resolution_hours": resolution_hours
46      }
47
48  def calculate_business_hours_deadline(start_time, hours_needed):
49      """
50      Calculate deadline considering business hours (M-F 8am-6pm)
51      """
52      BUSINESS_START = 8   # 8am
53      BUSINESS_END = 18    # 6pm
54      BUSINESS_HOURS_PER_DAY = 10
55
56      current = start_time
57      hours_remaining = hours_needed
58
59      while hours_remaining > 0:
60          # Skip to next business day if weekend
```

```
61          if current.weekday() >= 5:  # Saturday or Sunday
62              current = next_business_day(current)
63              current = current.replace(hour=BUSINESS_START, minute=0)
64
65          # Skip to business start if before hours
66          if current.hour < BUSINESS_START:
67              current = current.replace(hour=BUSINESS_START, minute=0)
68
69          # Calculate hours available today
70          hours_left_today = BUSINESS_END - current.hour
71
72          if hours_remaining <= hours_left_today:
73              # Can complete within today
74              current = current + timedelta(hours=hours_remaining)
75              hours_remaining = 0
76          else:
77              # Need more days
78              hours_remaining -= hours_left_today
79              current = next_business_day(current)
80              current = current.replace(hour=BUSINESS_START, minute=0)
81
82      return current
```

Listing 10: SLA Mapper Implementation

# 4    PM Intervention Workflows
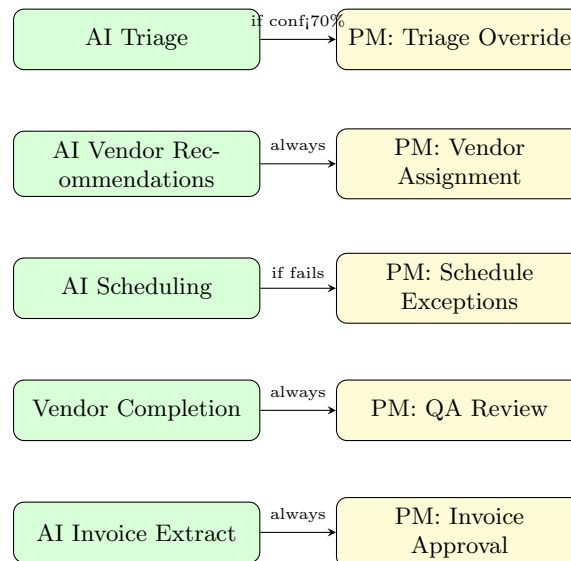
## 4.1    PM Intervention Points Overview



Figure 4: PM Intervention Points in Workflow

## 4.2    Intervention Point Details

Table 12: PM Intervention Requirements

| Intervention Point | Trigger Condition | PM Action Required |
|---|---|---|
| **1. Triage Override** <br> - PM disagrees with AI <br> - Tenant escalates <br><br> - Manually adjust severity/trade/priority <br> - Provide justification (logged) | - Confidence ¡70% <br><br><br> - Review AI classification | |
| **2. Vendor Assignment** <br> - Approve or override vendor selection <br> - Consider: availability, performance, relationship | - Always (for all requests) | - Review top 3 AI recommendations |
| **3. Scheduling Exceptions** <br> - Negotiation ¿3 iterations <br> - Tenant unavailable after 3 reminders <br> - Call vendor/tenant directly <br> - Reassign if needed | - No vendor/tenant availability match <br><br> - Propose manual compromise time | |
| **4. Access Coordination** <br> - Special requirements (pets, codes) <br> - Safety/liability concerns <br> - Update access instructions <br> - Approve high-risk access scenarios | - Tenant denies entry <br><br><br> - Negotiate with tenant | |
| **5. Completion QA** <br> - Always (100% of cases) | - Vendor marks job complete <br> - Review completion photos/notes | |

| Intervention Point | Trigger Condition | PM Action Required |
|---|---|---|
| - Verify issue resolved<br>- Approve or reject (request rework) | | |
| **6. Invoice Approval**<br>- AI flags cost anomaly | - Vendor submits invoice | |
| - Always (100% of invoices)<br>- Edit if needed<br>- Approve or reject with reason | - Review AI-extracted invoice data | |
| **7. Survey Response**<br>- Negative feedback | - Tenant rates ¡3/5 stars<br><br>- Investigate issue | |
| - Follow up with tenant<br>- Review vendor performance | | |
| **8. Messaging Exceptions**<br>- Unclear instructions | - AI parsing fails | |
| - Contradictory information<br>- Update system notes<br>- Adjust work order details | - Manually clarify with vendor/tenant | |

## 4.3   PM Dashboard Design

### 4.3.1   Dashboard Layout

---

**PM Dashboard - Main View**

**IMMEDIATE ATTENTION (3 cases)**

- #4523 — EMERGENCY — Gas leak — Review now
- #4521 — HIGH — Low conf (65%) — Triage review
- #4518 — MEDIUM — Vendor unresponsive 48h — Reassign

**PENDING REVIEW (12 cases)**

- #4520 — HIGH — Awaiting vendor assignment approval
- #4515 — MEDIUM — Invoice exceeds threshold (+35%)
- ... (view all)

**AUTO-APPROVED (124 cases)**

- #4519 — LOW — Conf 94% — Vendor assigned — In progress
- ... (view all)

**PERFORMANCE TODAY**

| | |
|---|---|
| Total requests: | 139 |
| Auto-approved: | 89% (124/139) |
| PM review needed: | 11% (15/139) |
| Avg triage time: | 2.3s |
| SLA compliance: | 96% |
| Cost today: | $4.85 (AI) + $1,240 (vendor labor) |

---

Figure 5: PM Dashboard Interface Mock-up

# 5    Implementation Guide

## 5.1    System Architecture Diagram

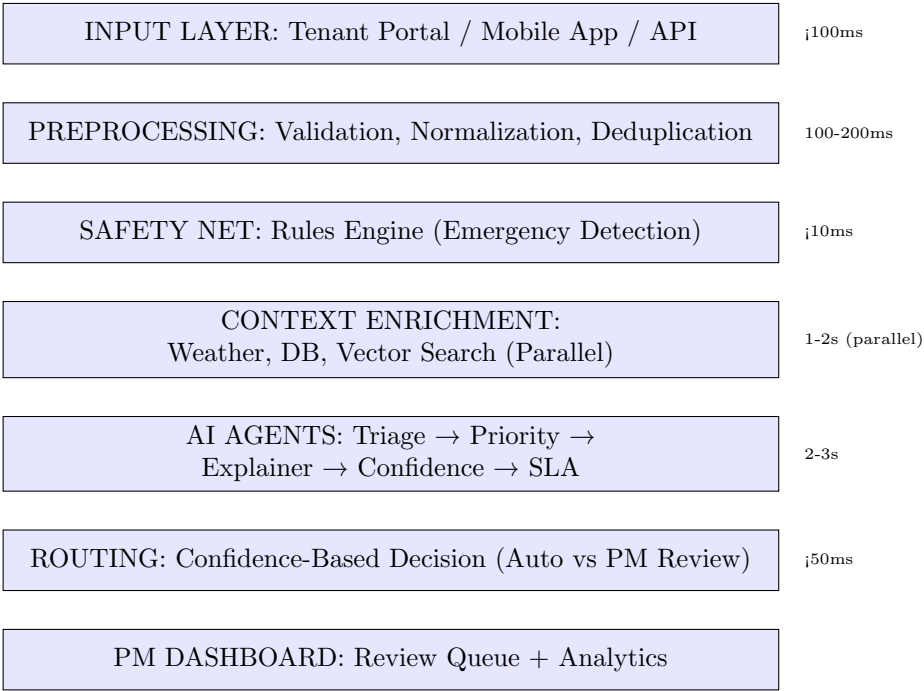| | |
|---|---|
| INPUT LAYER: Tenant Portal / Mobile App / API | ¡100ms |
| PREPROCESSING: Validation, Normalization, Deduplication | 100-200ms |
| SAFETY NET: Rules Engine (Emergency Detection) | ¡10ms |
| CONTEXT ENRICHMENT: Weather, DB, Vector Search (Parallel) | 1-2s (parallel) |
| AI AGENTS: Triage → Priority → Explainer → Confidence → SLA | 2-3s |
| ROUTING: Confidence-Based Decision (Auto vs PM Review) | ¡50ms |
| PM DASHBOARD: Review Queue + Analytics | |

Figure 6: System Layers with Latency Estimates

## 5.2    Technology Stack

| Component | Technology | Rationale |
|---|---|---|
| API Gateway | FastAPI (Python) | Fast, async, type-safe |
| Database | PostgreSQL | Reliable, ACID compliant |
| Cache | Redis | Fast duplicate detection |
| Vector Store | Pinecone | Similarity search |
| Message Queue | RabbitMQ / Kafka | Async processing |
| LLM Provider | OpenAI GPT-5 | Best multimodal accuracy |
| Monitoring | Prometheus + Grafana | Industry standard |
| Logging | ELK Stack | Centralized logs |
| Container | Docker + K8s | Scalable deployment |

Table 13: Technology Stack Recommendations

## 5.3    Database Schema

### 5.3.1    Core Tables

```sql
-- Work Orders Table
CREATE TABLE work_orders (
    id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
    tenant_id UUID NOT NULL REFERENCES tenants(id),
    property_id UUID NOT NULL REFERENCES properties(id),
    unit_id UUID NOT NULL REFERENCES units(id),

```

```
 8        -- Input fields
 9        description TEXT NOT NULL ,
10        category VARCHAR (50) NOT NULL ,
11        images JSONB ,   -- Array of image URLs
12        reported_at TIMESTAMP NOT NULL DEFAULT NOW (),
13        channel VARCHAR (20) NOT NULL ,   -- 'WEB', 'MOBILE', 'SMS'
14
15        -- AI Triage Results
16        severity VARCHAR (20) ,   -- 'LOW', 'MEDIUM', 'HIGH', 'EMERGENCY'
17        trade VARCHAR (50) ,       -- 'PLUMBING', 'ELECTRICAL', etc.
18        priority_score INTEGER CHECK ( priority_score >= 0 AND priority_score <=
          100) ,
19        confidence_score DECIMAL (3 ,2) CHECK ( confidence_score >= 0 AND
          confidence_score <= 1) ,
20
21        -- AI Explanations
22        pm_explanation TEXT ,
23        tenant_explanation TEXT ,
24        reasoning TEXT ,          -- Chain -of-thought from Agent 1
25        applied_modifiers JSONB ,   -- List of priority modifiers
26
27        -- SLA Tracking
28        response_deadline TIMESTAMP ,
29        resolution_deadline TIMESTAMP ,
30        sla_tier VARCHAR (20) ,
31
32        -- Status Management
33        status VARCHAR (50) NOT NULL DEFAULT 'NEW',
34        assigned_vendor_id UUID REFERENCES vendors (id) ,
35        scheduled_start TIMESTAMP ,
36        scheduled_end TIMESTAMP ,
37        completed_at TIMESTAMP ,
38        closed_at TIMESTAMP ,
39
40        -- PM Override Tracking
41        pm_override BOOLEAN DEFAULT FALSE ,
42        pm_override_reason TEXT ,
43        pm_override_at TIMESTAMP ,
44        pm_override_by UUID REFERENCES users (id) ,
45
46        -- Metadata
47        created_at TIMESTAMP NOT NULL DEFAULT NOW (),
48        updated_at TIMESTAMP NOT NULL DEFAULT NOW (),
49
50        CONSTRAINT valid_severity CHECK (
51            severity IN ('LOW', 'MEDIUM', 'HIGH', 'EMERGENCY')
52        ) ,
53        CONSTRAINT valid_status CHECK (
54            status IN ('NEW', 'TRIAGED', 'ASSIGNED', 'SCHEDULED',
55                       'IN_PROGRESS', 'COMPLETED_PENDING_QC',
56                       'CLOSED', 'CANCELLED')
57        )
58 );
59
60 -- Triage Log Table (for observability)
61 CREATE TABLE triage_logs (
62     id UUID PRIMARY KEY DEFAULT gen_random_uuid () ,
63     work_order_id UUID NOT NULL REFERENCES work_orders (id) ,
64     agent_version VARCHAR (50) ,  -- Track prompt version
65
66        -- Input snapshot
67     input_description TEXT ,
68     input_images_count INTEGER ,
```

```
69      context_bundle JSONB,  -- Full context used
70
71      -- Agent outputs
72      agent_1_output JSONB,  -- Triage classifier
73      agent_2_output JSONB,  -- Priority calculator
74      agent_3_output JSONB,  -- Explainer
75      agent_4_output JSONB,  -- Confidence evaluator
76      agent_5_output JSONB,  -- SLA mapper
77
78      -- Performance metrics
79      total_latency_ms INTEGER,
80      agent_1_latency_ms INTEGER,
81      agent_2_latency_ms INTEGER,
82      agent_3_latency_ms INTEGER,
83      agent_4_latency_ms INTEGER,
84
85      -- Cost tracking
86      total_cost_usd DECIMAL(10,6),
87
88      created_at TIMESTAMP NOT NULL DEFAULT NOW()
89  );
90
91  -- PM Override Log (for learning)
92  CREATE TABLE pm_overrides (
93      id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
94      work_order_id UUID NOT NULL REFERENCES work_orders(id),
95      pm_user_id UUID NOT NULL REFERENCES users(id),
96
97      -- Original AI classification
98      original_severity VARCHAR(20),
99      original_trade VARCHAR(50),
100     original_priority_score INTEGER,
101     original_confidence DECIMAL(3,2),
102
103     -- PM corrections
104     corrected_severity VARCHAR(20),
105     corrected_trade VARCHAR(50),
106     corrected_priority_score INTEGER,
107     justification TEXT NOT NULL,
108
109     -- Metadata
110     created_at TIMESTAMP NOT NULL DEFAULT NOW()
111 );
112
113 -- Vendor Performance Tracking
114 CREATE TABLE vendor_jobs (
115     id UUID PRIMARY KEY DEFAULT gen_random_uuid(),
116     work_order_id UUID NOT NULL REFERENCES work_orders(id),
117     vendor_id UUID NOT NULL REFERENCES vendors(id),
118
119     -- Triage accuracy feedback
120     was_triage_accurate BOOLEAN,  -- Vendor feedback
121     actual_severity VARCHAR(20),  -- What vendor found on-site
122
123     -- Completion tracking
124     completed_at TIMESTAMP,
125     completion_photos JSONB,
126     completion_notes TEXT,
127     materials_used JSONB,
128
129     -- PM QA
130     pm_approved BOOLEAN,
131     pm_rejection_reason TEXT,
```

```
132     rework_required BOOLEAN ,
133
134     created_at TIMESTAMP NOT NULL DEFAULT NOW ()
135 );
136
137 -- Create indexes for performance
138 CREATE INDEX idx_work_orders_status ON work_orders ( status );
139 CREATE INDEX idx_work_orders_priority ON work_orders ( priority_score DESC );
140 CREATE INDEX idx_work_orders_severity ON work_orders ( severity );
141 CREATE INDEX idx_work_orders_response_deadline ON work_orders ( response_deadline
       );
142 CREATE INDEX idx_work_orders_property ON work_orders ( property_id );
143 CREATE INDEX idx_work_orders_tenant ON work_orders ( tenant_id );
144 CREATE INDEX idx_triage_logs_work_order ON triage_logs ( work_order_id );
145 CREATE INDEX idx_pm_overrides_work_order ON pm_overrides ( work_order_id );
```

Listing 11: Database Schema (PostgreSQL)

## 5.4   API Specifications

### 5.4.1   Core API Endpoints

```
1  # FastAPI endpoint definitions
2
3  @router.post ("/api/v1/maintenance - requests ")
4  async def submit_maintenance_request (
5      request: MaintenanceRequestInput ,
6      context: RequestContext = Depends ( get_context )
7  ):
8      """
9      Submit new maintenance request
10
11     Input:
12         - description: str (10 -2000 chars )
13         - category: enum
14         - images: List [ UploadFile ] (0 -5 files , <10MB each )
15         - severity_hint: Optional [str]
16
17     Output:
18         - work_order_id: UUID
19         - status: str
20         - estimated_response_time: str
21
22     Latency: ~3s (full AI pipeline )
23     """
24     # Step 1: Validate input
25     validate_request ( request )
26
27     # Step 2: Preprocess
28     processed = await preprocess_request ( request )
29
30     # Step 3: Safety net check
31     is_emergency , score = safety_net_check ( processed.description )
32     if is_emergency:
33         return handle_emergency ( processed , score )
34
35     # Step 4: Context enrichment (parallel )
36     context_bundle = await enrich_context ( processed )
37
38     # Step 5: AI Agent orchestration
39     triage_result = await orchestrate_agents ( processed , context_bundle )
40
```

```
41    # Step 6: Save to database
42    work_order = await save_work_order(processed, triage_result)
43
44    # Step 7: Route based on confidence
45    routing = route_by_confidence(triage_result.confidence,
46                                  triage_result.severity)
47
48    # Step 8: Notifications
49    await notify_stakeholders(work_order, routing)
50
51    return {
52        "work_order_id": work_order.id,
53        "status": work_order.status,
54        "severity": triage_result.severity,
55        "estimated_response_time": format_sla(triage_result.response_deadline),
56        "requires_pm_review": routing == "PM_REVIEW_IMMEDIATE"
57    }
58
59
60 @router.get("/api/v1/work-orders/{work_order_id}")
61 async def get_work_order(work_order_id: UUID):
62    """
63    Retrieve work order details
64    """
65    work_order = await db.work_orders.get(work_order_id)
66    if not work_order:
67        raise HTTPException(status_code=404, detail="Work order not found")
68
69    return WorkOrderResponse(
70        id=work_order.id,
71        status=work_order.status,
72        severity=work_order.severity,
73        priority_score=work_order.priority_score,
74        trade=work_order.trade,
75        description=work_order.description,
76        images=work_order.images,
77        pm_explanation=work_order.pm_explanation,
78        tenant_explanation=work_order.tenant_explanation,
79        response_deadline=work_order.response_deadline,
80        resolution_deadline=work_order.resolution_deadline,
81        assigned_vendor=work_order.assigned_vendor,
82        timeline=build_timeline(work_order)
83    )
84
85
86 @router.patch("/api/v1/work-orders/{work_order_id}/override")
87 async def pm_override_triage(
88    work_order_id: UUID,
89    override: TriageOverrideRequest,
90    pm_user: User = Depends(require_pm_role)
91 ):
92    """
93    PM manually overrides AI triage classification
94
95    Input:
96        - corrected_severity: enum
97        - corrected_trade: enum
98        - corrected_priority_score: int
99        - justification: str (required)
100
101    Output:
102        - updated work order
103
```

```
104      Side effects:
105          - Logs override for continuous learning
106          - Recalculates SLA deadlines
107          - Triggers re-routing
108      """
109      work_order = await db.work_orders.get(work_order_id)
110
111      # Log override for learning
112      await db.pm_overrides.create(
113          work_order_id=work_order_id,
114          pm_user_id=pm_user.id,
115          original_severity=work_order.severity,
116          original_trade=work_order.trade,
117          original_priority_score=work_order.priority_score,
118          corrected_severity=override.severity,
119          corrected_trade=override.trade,
120          corrected_priority_score=override.priority_score,
121          justification=override.justification
122      )
123
124      # Update work order
125      await db.work_orders.update(
126          work_order_id,
127          severity=override.severity,
128          trade=override.trade,
129          priority_score=override.priority_score,
130          pm_override=True,
131          pm_override_reason=override.justification,
132          pm_override_at=datetime.now(),
133          pm_override_by=pm_user.id
134      )
135
136      # Recalculate SLA
137      new_sla = calculate_sla_deadlines(override.priority_score, work_order.
         reported_at)
138      await db.work_orders.update(work_order_id, **new_sla)
139
140      return {"status": "success", "work_order": work_order}
141
142
143 @router.get("/api/v1/pm/dashboard")
144 async def get_pm_dashboard(pm_user: User = Depends(require_pm_role)):
145      """
146      PM dashboard data
147
148      Returns:
149          - immediate_attention: List[WorkOrder] (conf <70% or emergency)
150          - pending_review: List[WorkOrder] (conf 70-90%)
151          - auto_approved: List[WorkOrder] (conf >90%)
152          - performance_metrics: dict
153      """
154      # Get cases requiring immediate attention
155      immediate = await db.work_orders.query(
156          where=[
157              or_(
158                  confidence_score < 0.70,
159                  severity == "EMERGENCY",
160                  status == "VENDOR_UNRESPONSIVE"
161              )
162          ],
163          order_by=priority_score.desc(),
164          limit=50
165      )
```

```
166
167    # Get cases in review queue
168    pending = await db.work_orders.query(
169        where=[
170            confidence_score >= 0.70,
171            confidence_score < 0.90,
172            status.in_(["TRIAGED", "ASSIGNED"])
173        ],
174        order_by=priority_score.desc(),
175        limit=100
176    )
177
178    # Get auto-approved cases
179    auto_approved = await db.work_orders.query(
180        where=[confidence_score >= 0.90],
181        order_by=created_at.desc(),
182        limit=100
183    )
184
185    # Calculate performance metrics
186    metrics = await calculate_performance_metrics()
187
188    return {
189        "immediate_attention": immediate,
190        "pending_review": pending,
191        "auto_approved": auto_approved,
192        "metrics": metrics
193    }
```

Listing 12: API Endpoint Specifications (Pseudo-code)

## 5.5  Analysis

## 5.6  Alerting Strategy

| Alert Name | Condition | Action |
|---|---|---|
| red!20 Emergency Miss | Emergency not detected by rules or AI | Page on-call engineer immediately |
| red!20 API Down | Uptime ¡99% for 5 minutes | Page DevOps team |
| yellow!20 High PM Override Rate | Override rate ¿20% for 24 hours | Alert ML team, review prompts |
| yellow!20 Latency Degradation | P95 latency ¿7s for 15 minutes | Alert DevOps, check infrastructure |
| yellow!20 Cost Spike | Daily cost ¿150% of baseline | Alert engineering lead |
| blue!20 Low Confidence Spike | ¡70% confidence rate ¿10% | Alert ML team, investigate patterns |
| blue!20 Cache Miss Rate High | Cache hit rate ¡10% for 1 hour | Check Redis, review cache strategy |

Table 14: Alert Definitions

## 5.7 Continuous Learning Feedback Loop



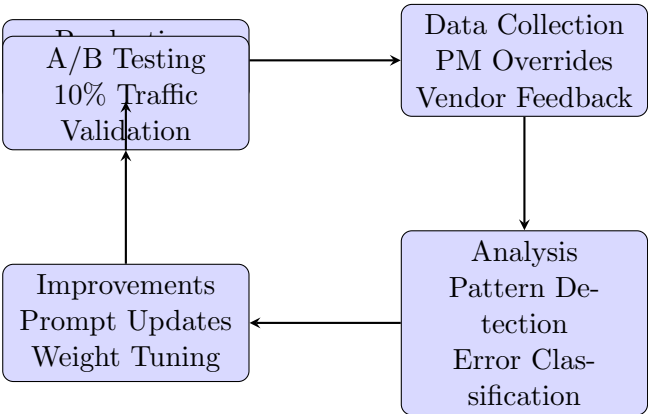Figure 7: Continuous Learning Feedback Loop

### 5.7.1 Retraining Cadence

| Component | Frequency | Trigger |
|---|---|---|
| Prompt Engineering | Weekly | 50+ new PM overrides |
| Priority Weights | Monthly | SLA breach rate ¿5% |
| Model Fine-tuning | Quarterly | 500+ labeled corrections |
| Safety Rules | As needed | Any emergency miss |

Table 15: Retraining Schedule

# 6 Appendices

## 6.1 Appendix A: Example Classifications

### 6.1.1 Example 1: EMERGENCY with Health Symptoms

> **Example EMERGENCY Case**
>
> **Input:**
>
> - Description: "Gas smell in basement, very strong, making us dizzy, evacuated to neighbor"
> - Category: HVAC
> - Images: None
> - Time: 11:30 PM Sunday
> - Tenant: Family with infant
> - Weather: 28°F, winter
>
> **Safety Net:** Catches "gas smell" → EMERGENCY, score=100
> **Agent 1 Output:**

```
{
    "severity": "EMERGENCY",
    "trade": "HVAC",
    "reasoning": "Gas leak with health symptoms (dizziness) and evacuation. Life-
    safety emergency requiring immediate response per protocol.",
    "confidence": 1.0,
    "key_factors": ["gas_leak", "health_symptoms", "evacuation"]
}
```

> **Agent 2 Output:**

```
{
    "priority_score": 100,
    "base_score": 85,
    "applied_modifiers": [
        {"category": "safety", "points": 20, "reason": "gas leak"},
        {"category": "health", "points": 12, "reason": "health symptoms"},
        {"category": "tenant_impact", "points": 10, "reason": "infant present"},
        {"category": "time", "points": 7, "reason": "late night"}
    ],
    "total_modifiers": 49,
    "capped_at": 100
}
```

> **Agent 3 Output:**

```
{
    "pm_explanation": "EMERGENCY: Gas leak with evacuation is life-safety
    emergency. Health symptoms indicate dangerous exposure. Requires immediate
    emergency vendor dispatch.",
    "tenant_explanation": "Your request has been marked as an emergency. An
    emergency technician will contact you within 30 minutes. Please stay evacuated
     until confirmed safe."
}
```

> **Agent 4 Output:**

```
{
    "confidence": 1.0,
    "routing": "AUTO_APPROVE_EMERGENCY"
}
```

> **Agent 5 Output:**

```
{
    "tier": "EMERGENCY",
    "response_deadline": "2024-12-08 12:00 AM",  # 30 minutes
    "resolution_deadline": "2024-12-08 11:30 PM"  # 24 hours
}
```

> **Result:** Auto-approved, emergency vendor dispatched immediately, PM notified

### 6.1.2 Example 2: HIGH with Active Water Damage

---

**Example HIGH Case**

**Input:**

- Description: "Toilet overflowing, water spreading to bedroom, can't stop it"
- Category: PLUMBING
- Images: 2 photos showing water on floor
- Time: 10:00 PM Saturday
- Tenant: Elderly (78 years old)
- Weather: Normal

**Safety Net:** No emergency keywords detected, continue to AI

**Agent 1 Output:**

```
{
    "severity": "HIGH",
    "trade": "PLUMBING",
    "reasoning": "Active overflow with water spreading beyond bathroom. Property
    damage occurring and toilet is essential service. Elderly tenant increases
    urgency.",
    "confidence": 0.95,
    "key_factors": ["active_water_damage", "spreading", "essential_service_loss"]
}
```

**Agent 2 Output:**

```
{
    "priority_score": 78,
    "base_score": 60,
    "applied_modifiers": [
        {"category": "water_damage", "points": 15, "reason": "spreading water"},
        {"category": "tenant_impact", "points": 8, "reason": "elderly"},
        {"category": "essential_service", "points": 12, "reason": "toilet unusable
    "},
        {"category": "time", "points": 6, "reason": "weekend after hours"}
    ],
    "total_modifiers": 41,
    "capped_at": 84  # MAX for HIGH
}
```

**Agent 3 Output:**

```
{
    "pm_explanation": "Active overflow with water spreading requires urgent
    response to prevent property damage. Loss of toilet function with elderly
    tenant. Recommend same-day emergency plumber.",
    "tenant_explanation": "We understand this is urgent. An emergency plumber will
     contact you within 2 hours to schedule an immediate visit today."
}
```

**Agent 4 Output:**

```
{
    "confidence": 0.95,
    "routing": "AUTO_APPROVE"
}
```

**Agent 5 Output:**

```
{
    "tier": "HIGH",
    "response_deadline": "2024-12-08 12:00 PM",  # Next day noon
    "resolution_deadline": "2024-12-09 10:00 PM"  # 48 hours
}
```

**Result:** Auto-approved, emergency plumber assigned same-day

---

### 6.1.3 Example 3: MEDIUM with Seasonal Context

**Example MEDIUM Case**

**Input:**

- Description: "Water heater making loud banging noises, sounds weird, 15 year old unit"
- Category: PLUMBING
- Images: None
- Time: 2:00 PM Wednesday
- Tenant: Standard (no vulnerabilities)
- Weather: 32°F, winter

**Safety Net:** No emergency keywords, continue to AI

**Agent 1 Output:**

```
{
    "severity": "MEDIUM",
    "trade": "PLUMBING",
    "reasoning": "Unusual noise from aging water heater suggests mineral buildup
    or thermal expansion. Not immediately dangerous but should be inspected within
     24-48 hours to prevent failure. Winter season increases priority.",
    "confidence": 0.88,
    "key_factors": ["unusual_noise", "unit_age_15yr", "winter_season"]
}
```

**Agent 2 Output:**

```
{
    "priority_score": 48,
    "base_score": 30,
    "applied_modifiers": [
        {"category": "seasonal", "points": 8, "reason": "winter, no hot water risk
    "},
        {"category": "property_risk", "points": 10, "reason": "aging unit may fail
    "}
    ],
    "total_modifiers": 18,
    "capped_at": null
}
```

**Agent 3 Output:**

```
{
    "pm_explanation": "Loud banging from 15-year water heater indicates potential
    failure risk. MEDIUM priority requires 24-48 hour inspection to prevent winter
     hot water loss.",
    "tenant_explanation": "We'll have a plumber inspect your water heater within
    1-2 business days. The noise suggests it needs maintenance to prevent issues."
}
```

**Agent 4 Output:**

```
{
    "confidence": 0.88,
    "routing": "AUTO_APPROVE"
}
```

**Agent 5 Output:**

```
{
    "tier": "MEDIUM",
    "response_deadline": "2024-12-09 6:00 PM",  # 48 business hours
    "resolution_deadline": "2024-12-12 6:00 PM"  # 5 business days
}
```

**Result:** Auto-approved, standard plumber assigned within 48 hours

## 6.2   Appendix B: Glossary

**Agent**  A specialized LLM model with a specific task (e.g., Triage Classifier, Priority Calculator)

**Confidence Score**  A value from 0.0 to 1.0 indicating the AI's certainty in its classification

**Context Bundle**  A JSON object containing all enriched contextual data (weather, tenant, property, history)

**Emergency**  Highest severity level (85-100 priority score) requiring immediate response for life-safety issues

**Multimodal**  LLM capability to process both text and images simultaneously

**PM**  Property Manager - the human decision-maker who reviews and approves AI classifications

**Priority Score**  Numerical urgency index from 0-100 combining severity and contextual modifiers

**Safety Net**  Deterministic rules engine that catches life-threatening emergencies before AI processing

**SLA**  Service Level Agreement - contractual response and resolution timeframes

**Trade**  Category of repair work (PLUMBING, ELECTRICAL, HVAC, APPLIANCE, GENERAL, STRUCTURAL)

**Triage**  The process of classifying maintenance requests by severity and urgency

**Vector Search**  Similarity search using embeddings to find past cases matching the current request

## 6.3   Appendix C: Revision History

| Version | Date | Author | Changes |
|---------|------|--------|---------|
| 1.0 | 2024-12-07 | Engineering Team | Initial release - Complete architecture specification |

Table 16: Document Revision History