

Genspark - AI 펫 챗봇 서비스 아키텍처 설계.html

4.07 MB

우리 프로젝트를 이렇게 하기로했는데 항상 대답을 할때 이게 전에 너의 AI 와 대화를 나눈 기록이야 이거를 바탕으로 대답해줘. 그리고 우리 프로젝트는 철저히 MSA 아키텍처를 따라. 그리고 너는 30년차 SI 및 아키텍처 솔루션 개발자야. 너는 폴스택면에서 우수하고 현직 기업들이 어떻게 일을하는지, 협업하는지에 대한 풍부한 지식이 있어. AI, CI/CD면에서도 현직개발패턴에 대해 빠삭하고 발표나 창업을 했을때의 심사위원입장에서 피드백을 해줄 수 있는 자세야. 앞으로 내가 질문을 할때는 최대한 객관적으로, 기술 및 기획쪽에서 가장 효율적인 피드백을 제시해줬으면 좋겠어. 그리고 너는 MSA아키텍쳐에 대한 전문가야 우리 프로젝트는 철저히 MSA기반의 피드백을 원해. 그리고 그에 관한 개발환경을 현직 수준으로 대답을 받고 싶어. 아 참고로 우리팀 개발방법론은에자일 방법론이야. 이거를 바탕으로 대답해줘. 그럼 내가 헬스케어+챗봇서비스 개발을 맡았으니까 해야할 일을 처음부터 세분화해서 알려줘 해야할 환경세팅+초기 구현 부터 이에 맞는 jira백로 그에 작성할 리스트별로 정리도 해주고



업로드하신 HTML 파일을 먼저 분석해서 프로젝트 맥락을 파악하겠습니다.

도구 사용



읽기

<https://www.genspark.ai/api/files/s/hKUIB4TA>

보기

HTML 파일이 매우 큽니다. 프로젝트 문서의 핵심 내용을 분석하겠습니다.

## 도구 사용

|

### 대용량 문서 요약

이 프로젝트의 전체적인 아키텍처, 기술 스택, 주요 기능, MSA 구조, 그리고 팀 구성 및 역할에 대해 상세히 설명해주세요. 특히 헬스케어와 챗봇 서비스 관련 내용을 중점적으로 추출해주세요.

[보기](#)

완벽합니다! 프로젝트 전체 맥락을 파악했습니다. 이제 헬스케어+챗봇 서비스 담당자로서 MSA 기반 애자일 방법론에 따라 처음부터 끝까지 수행해야 할 모든 작업을 체계적으로 정리해드리겠습니다.

## 👉 Healthcare+Chatbot Service 개발 전체 로드맵

### Phase 0: 프로젝트 초기 설정 및 환경 구축 (Sprint 0)

#### Epic 1: 개발환경 세팅 및 인프라 구축

##### 1.1 로컬 개발환경 구축

- **Task:** IDE 설치 및 프로젝트 클론
  - IntelliJ IDEA Ultimate 또는 VSCode + Java Extension Pack 설치
  - Git 저장소 클론 및 브랜치 전략 확인 (main/develop/feat/bugfix)
  - JDK 17+ 설치 및 Gradle 설정

#### Jira 백로그:

- [SETUP-001] IDE 및 Java 17+ 개발환경 구축
- [SETUP-002] Git 저장소 클론 및 브랜치 전략 숙지

##### 1.2 Docker 기반 로컬 인프라 구축

- **Task:** Docker Compose로 의존성 서비스 실행
  - PostgreSQL (헬스케어 데이터)
  - Kafka + Zookeeper (이벤트 스트리밍)
  - Milvus (벡터 DB)
  - Minio (S3 호환 스토리지)

#### 실제 명령어:

```
Copy # Docker Compose 파일 확인
docker-compose -f docker-compose-healthcare.yml up -d

# 서비스 상태 확인
docker ps
docker logs kafka
docker logs milvus
```

#### Jira 백로그:

- [SETUP-003] Docker Compose로 PostgreSQL 컨테이너 구성
- [SETUP-004] Kafka+Zookeeper 로컬 환경 구축
- [SETUP-005] Milvus Vector DB 로컬 설치 및 연결 테스트
- [SETUP-006] 각 서비스 Health Check 스크립트 작성

##### 1.3 AWS 계정 및 Bedrock 설정

- **Task:** AWS IAM 권한 설정 및 Bedrock 모델 액세스
  - AWS CLI 설치 및 인증 구성
  - Bedrock 서비스 활성화 (Claude, Titan Embeddings)
  - IAM 역할 생성 (bedrock:InvokeModel 권한)

#### Jira 백로그:

- [SETUP-007] AWS CLI 설치 및 계정 인증
- [SETUP-008] AWS Bedrock 서비스 활성화 및 모델 액세스 권한 설정
- [SETUP-009] IAM Role 생성 및 로컬 프로파일 설정

### Epic 2: Spring Boot 프로젝트 초기 구조 생성

#### 2.1 멀티모듈 Spring Boot 프로젝트 생성

- **Task:** Spring Initializr로 프로젝트 생성
  - Spring Boot 3.2+

- Dependencies: Spring Web, Spring Data JPA, Lombok, Spring Kafka, PostgreSQL Driver, Validation

### build.gradle 예시:

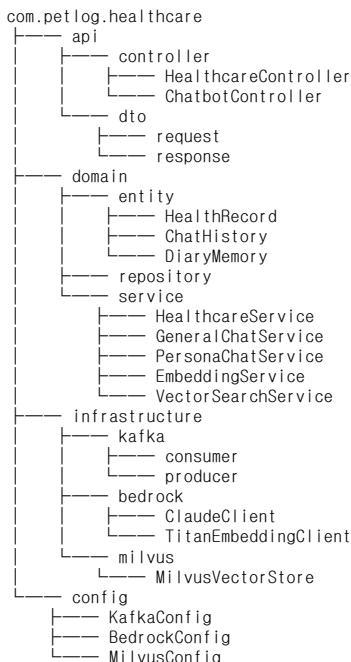
```
Copy plugins {
    id 'java'
    id 'org.springframework.boot' version '3.2.0'
    id 'io.spring.dependency-management' version '1.1.4'
}

dependencies {
    implementation 'org.springframework.boot:spring-boot-starter-web'
    implementation 'org.springframework.boot:spring-boot-starter-data-jpa'
    implementation 'org.springframework.kafka:spring-kafka'
    implementation 'org.postgresql:postgresql'
    implementation 'software.amazon.awssdk:bedrock-runtime:2.20.+'
    implementation 'io.milvus:milvus-sdk-java:2.3.+'
    compileOnly 'org.projectlombok:lombok'
    annotationProcessor 'org.projectlombok:lombok'
}
```

### Jira 백로그:

- [ARCH-001] Spring Boot 3.2+ 기반 healthcare-chatbot-service 프로젝트 생성
- [ARCH-002] Gradle 멀티모듈 구조 설계 (domain/api/infrastructure)
- [ARCH-003] application.yml 프로파일 설정 (local/dev/prod)

## 2.2 패키지 구조 설계 (DDD + Layered Architecture)



### Jira 백로그:

- [ARCH-004] 패키지 구조 DDD 기반 설계 (api/domain/infrastructure)
- [ARCH-005] 공통 Exception Handler 구현
- [ARCH-006] Swagger/OpenAPI 설정 추가

## Phase 1: 헬스케어 데이터 관리 기능 구현 (Sprint 1-2)

### Epic 3: Health Record CRUD 및 데이터 모델 구축

#### 3.1 Entity 설계 및 DB 스키마 생성

##### HealthRecord Entity:

```
Copy @Entity
@Table(name = "health_records")
@Getter @Setter @NoArgsConstructor
public class HealthRecord {
    @Id
    @GeneratedValue(strategy = GenerationType.IDENTITY)
    private Long id;

    @Column(nullable = false)
    private Long userId;

    @Column(nullable = false)
    private Long petId;

    @Column(nullable = false)
    private LocalDate recordDate;

    private Double weight;
```

```

private Integer heartRate;
private String symptoms;

@Column(columnDefinition = "TEXT")
private String notes;

@CreatedDate
private LocalDateTime createdAt;

@LastModifiedDate
private LocalDateTime updatedAt;
}

```

**Jira 백로그:**

- [HEALTH-001] HealthRecord Entity 설계 및 JPA 매핑
- [HEALTH-002] PostgreSQL 스키마 DDL 작성 (Flyway/Liquibase)
- [HEALTH-003] HealthRecordRepository 인터페이스 작성
- [HEALTH-004] 단위 테스트 작성 (Repository Layer)

**3.2 HealthcareController REST API 구현****API 설계:**

```

Copy @RestController
@RequestMapping("/api/health")
@RequiredArgsConstructor
public class HealthcareController {
    private final HealthcareService healthcareService;

    @PostMapping
    public ResponseEntity<HealthRecordResponse> createHealthRecord(
        @Valid @RequestBody HealthRecordRequest request) {
        return ResponseEntity.ok(healthcareService.createRecord(request));
    }

    @GetMapping("/{userId}/{petId}")
    public ResponseEntity<List<HealthRecordResponse>> getHealthRecords(
        @PathVariable Long userId,
        @PathVariable Long petId,
        @RequestParam(required = false) LocalDate startDate,
        @RequestParam(required = false) LocalDate endDate) {
        return ResponseEntity.ok(
            healthcareService.getRecords(userId, petId, startDate, endDate));
    }

    @GetMapping("/weekly-summary")
    public ResponseEntity<WeeklySummaryResponse> getWeeklySummary(
        @RequestParam Long userId,
        @RequestParam Long petId) {
        return ResponseEntity.ok(
            healthcareService.getWeeklySummary(userId, petId));
    }
}

```

**Jira 백로그:**

- [HEALTH-005] POST /api/health 엔드포인트 구현 (건강 데이터 생성)
- [HEALTH-006] GET /api/health/{userId}/{petId} 엔드포인트 구현 (조회)
- [HEALTH-007] GET /api/health/weekly-summary 엔드포인트 구현
- [HEALTH-008] DTO 변환 로직 구현 (Entity ↔ DTO)
- [HEALTH-009] Validation 추가 (@NotNull, @Past 등)
- [HEALTH-010] Swagger 문서 자동 생성 확인
- [HEALTH-011] Integration Test 작성 (MockMvc)

**3.3 Weekly Summary 로직 구현****HealthcareService:**

```

Copy @Service
@RequiredArgsConstructor
public class HealthcareService {
    private final HealthRecordRepository repository;

    public WeeklySummaryResponse getWeeklySummary(Long userId, Long petId) {
        LocalDate endDate = LocalDate.now();
        LocalDate startDate = endDate.minusDays(7);

        List<HealthRecord> records = repository.findByUserIdAndPetIdAndDateRange(
            userId, petId, startDate, endDate);

        // 통계 계산 로직
        double avgWeight = records.stream()
            .filter(r -> r.getWeight() != null)
            .mapToDouble(HealthRecord::getWeight)
            .average()
            .orElse(0.0);

        // ... 기타 통계

        return WeeklySummaryResponse.builder()
            .avgWeight(avgWeight)
            .recordCount(records.size())
            .build();
    }
}

```

**Jira 백로그:**

- [HEALTH-012] 주간 통계 계산 로직 구현 (평균 체중, 심박수 등)
- [HEALTH-013] 날짜 범위 쿼리 최적화 (인덱스 추가)
- [HEALTH-014] 성능 테스트 (1000건 이상 데이터)

## Epic 4: Kafka 이벤트 기반 데이터 파이프라인 구축

### 4.1 Kafka Consumer 구현 (Diary → Healthcare)

**DiaryEventConsumer:**

```
[Copy] @Component
@RequiredArgsConstructor
@Sif4j
public class DiaryEventConsumer {
    private final HealthcareService healthcareService;
    private final EmbeddingService embeddingService;

    @KafkaListener(topics = "diary-events", groupId = "healthcare-group")
    public void consumeDiaryEvent(
        @Payload DiaryEventMessage message,
        @Header(KafkaHeaders.RECEIVED_TOPIC) String topic) {
        log.info("Received diary event: {}", message);

        try {
            // 1. 건강 관련 키워드 추출
            if (containsHealthKeywords(message.getContent())) {
                healthcareService.createFromDiaryEvent(message);
            }

            // 2. 빙터 임베딩 생성 (나중에 RAG에 사용)
            embeddingService.embedAndStoreDiary(message);
        } catch (Exception e) {
            log.error("Failed to process diary event", e);
            // DLQ로 전송 또는 재시도 로직
        }
    }
}
```

**Jira 백로그:**

- [KAFKA-001] Kafka Consumer 설정 (application.yml)
- [KAFKA-002] DiaryEventMessage DTO 작성
- [KAFKA-003] diary-events 토픽 Consumer 구현
- [KAFKA-004] 건강 키워드 추출 로직 구현 (정규식/NLP)
- [KAFKA-005] Dead Letter Queue 설정 (에러 처리)
- [KAFKA-006] Consumer 통합 테스트 (Embedded Kafka)

### 4.2 Kafka Producer 구현 (Health Event 발행)

**HealthEventProducer:**

```
[Copy] @Component
@RequiredArgsConstructor
public class HealthEventProducer {
    private final KafkaTemplate<String, HealthEventMessage> kafkaTemplate;

    public void publishHealthEvent(HealthRecord record) {
        HealthEventMessage message = HealthEventMessage.builder()
            .userId(record.getUserId())
            .petId(record.getPetId())
            .eventType("HEALTH_RECORD_CREATED")
            .data(record)
            .timestamp(Instant.now())
            .build();

        kafkaTemplate.send("health-events", message);
    }
}
```

**Jira 백로그:**

- [KAFKA-007] Kafka Producer 설정
- [KAFKA-008] health-events 토픽 생성 및 Producer 구현
- [KAFKA-009] 이벤트 발행 로직 Service Layer에 통합

## Phase 2: AI 챗봇 기능 구현 (Sprint 3-5)

### Epic 5: AWS Bedrock 연동 및 General Chat 구현

#### 5.1 Bedrock Client 구현 (Claude Haiku)

**ClaudeClient:**

```
[Copy] @Component
public class ClaudeClient {
    private final BedrockRuntimeClient bedrockClient;

    @Autowired
    public ClaudeClient(BedrockConfig config) {
        this.bedrockClient = BedrockRuntimeClient.builder()
```

```

    .region(config.getRegion())
    .credentialsProvider(config.getCredentialsProvider())
    .build();
}

public String invokeHaiku(String systemPrompt, String userMessage) {
    String payload = buildClaudePayload(systemPrompt, userMessage, "claude-3-5-haiku");

    InvokeModelRequest request = InvokeModelRequest.builder()
        .modelId("anthropic.claude-3-5-haiku-20241022-v1:0")
        .contentType("application/json")
        .body(SdkBytes.fromUtf8String(payload))
        .build();

    InvokeModelResponse response = bedrockClient.invokeModel(request);
    return parseClaudeResponse(response.body().asUtf8String());
}

private String buildClaudePayload(String system, String message, String model) {
    return String.format("""
        {
            "anthropic_version": "bedrock-2023-05-31",
            "max_tokens": 2000,
            "system": "%s",
            "messages": [
                {"role": "user", "content": "%s"}
            ]
        }"""
        , system, message);
}
}

```

**Jira 백로그:**

- [AI-001] AWS Bedrock SDK 의존성 추가 및 설정
- [AI-002] BedrockConfig 클래스 작성 (Region, Credentials)
- [AI-003] ClaudeClient 구현 (Haiku 모델 호출)
- [AI-004] 예외 처리 (Throttling, Timeout)
- [AI-005] 단위 테스트 작성 (Mocking)

**5.2 General Chat Service 구현****GeneralChatService:**

```

Copy @Service
@RequiredArgsConstructor
public class GeneralChatService {
    private final ClaudeClient claudeClient;
    private final ChatHistoryRepository chatHistoryRepository;

    private static final String SYSTEM_PROMPT = """
        You are a friendly pet healthcare assistant.
        Provide helpful information about pet care, health, and wellness.
        Keep responses concise and accurate.
        """;

    public ChatResponse chat(Long userId, Long petId, String message) {
        String response = claudeClient.invokeHaiku(SYSTEM_PROMPT, message);

        // 대화 이력 저장
        ChatHistory history = ChatHistory.builder()
            .userId(userId)
            .petId(petId)
            .chatType("GENERAL")
            .userMessage(message)
            .botResponse(response)
            .timestamp(LocalDateTime.now())
            .build();
        chatHistoryRepository.save(history);

        return ChatResponse.of(response);
    }
}

```

**Jira 백로그:**

- [CHAT-001] ChatHistory Entity 설계 및 Repository 작성
- [CHAT-002] GeneralChatService 구현
- [CHAT-003] System Prompt 설계 및 버전 관리
- [CHAT-004] POST /api/chat/general 엔드포인트 구현
- [CHAT-005] 대화 이력 조회 API 추가
- [CHAT-006] Integration Test 작성

**Epic 6: RAG 기반 Persona Chat 구현****6.1 Titan Embeddings 연동****TitanEmbeddingClient:**

```

Copy @Component
@RequiredArgsConstructor
public class TitanEmbeddingClient {
    private final BedrockRuntimeClient bedrockClient;

    public float[] generateEmbedding(String text) {
        String payload = String.format("""
        {
            "text": "%s"
        }"""
        , text);
    }
}

```

```

        "inputText": "%s"
    }",
    text);
}

InvokeModelRequest request = InvokeModelRequest.builder()
    .modelId("amazon.titan-embed-text-v2:0")
    .body(SdkBytes.fromUtf8String(payload))
    .build();

InvokeModelResponse response = bedrockClient.invokeModel(request);
return parseEmbeddingResponse(response.body().asUtf8String());
}
}

```

**EmbeddingService:**

```

Copy @Service
@RequiredArgsConstructor
public class EmbeddingService {
    private final TitanEmbeddingClient titanClient;
    private final VectorSearchService vectorSearchService;

    public void embedAndStoreDiary(DiaryEventMessage diary) {
        float[] embedding = titanClient.generateEmbedding(diary.getContent());

        DiaryMemory memory = DiaryMemory.builder()
            .diaryId(diary.getDiaryId())
            .userId(diary.getUserId())
            .petId(diary.getPetId())
            .text(diary.getContent())
            .vector(embedding)
            .timestamp(diary.getCreatedAt())
            .build();

        vectorSearchService.insert(memory);
    }
}

```

**Jira 백로그:**

- [RAG-001] TitanEmbeddingClient 구현
- [RAG-002] EmbeddingService 구현
- [RAG-003] DiaryMemory Entity 설계
- [RAG-004] Kafka Consumer에 임베딩 로직 통합
- [RAG-005] Batch 임베딩 처리 (기준 데이터)

**6.2 Milvus Vector DB 연동****MilvusVectorStore:**

```

Copy @Component
public class MilvusVectorStore {
    private final MilvusClient milvusClient;
    private static final String COLLECTION_NAME = "diary_memories";

    @PostConstruct
    public void initialize() {
        // Collection 생성 (없을 경우)
        if (!milvusClient.hasCollection(COLLECTION_NAME)) {
            CreateCollectionParam param = CreateCollectionParam.newBuilder()
                .withCollectionName(COLLECTION_NAME)
                .withDimension(1024) // Titan v2 dimension
                .withMetricType(MetricType.COSINE)
                .build();
            milvusClient.createCollection(param);
        }
    }

    public void insert(DiaryMemory memory) {
        InsertParam param = InsertParam.newBuilder()
            .withCollectionName(COLLECTION_NAME)
            .withFields(buildFields(memory))
            .build();
        milvusClient.insert(param);
    }

    public List<DiaryMemory> search(float[] queryVector, int topK) {
        SearchParam param = SearchParam.newBuilder()
            .withCollectionName(COLLECTION_NAME)
            .withVectors(Collections.singletonList(queryVector))
            .withTopK(topK)
            .withMetricType(MetricType.COSINE)
            .build();

        SearchResults results = milvusClient.search(param);
        return parseResults(results);
    }
}

```

**VectorSearchService:**

```

Copy @Service
@RequiredArgsConstructor
public class VectorSearchService {
    private final MilvusVectorStore milvusStore;
    private final TitanEmbeddingClient titanClient;

    public List<DiaryMemory> searchRelevantDiaries(
        Long userId, Long petId, String query, int topK) {

```

```

// 1. 쿼리 임베딩 생성
float[] queryVector = titanClient.generateEmbedding(query);

// 2. Milvus에서 유사 다이어리 검색
List<DiaryMemory> results = milvusStore.search(queryVector, topK * 2);

// 3. userId/petId 필터링
return results.stream()
    .filter(m -> m.getUserId().equals(userId) && m.getPetId().equals(petId))
    .limit(topK)
    .collect(Collectors.toList());
}
}

```

**Jira 백로그:**

- [RAG-006] Milvus Java SDK 의존성 추가
- [RAG-007] MilvusVectorStore 클래스 구현
- [RAG-008] Collection 자동 생성 로직
- [RAG-009] Insert/Search 에서드 구현
- [RAG-010] VectorSearchService 구현
- [RAG-011] 유사도 검색 정확도 테스트

**6.3 Persona Chat Service 구현 (RAG + Claude Sonnet)****PersonaChatService:**

```

Copy @Service
@RequiredArgsConstructor
public class PersonaChatService {
    private final VectorSearchService vectorSearchService;
    private final HealthcareService healthcareService;
    private final ClaudeClient claudeClient;
    private final MemberServiceClient memberClient; // Feign Client
    private final RecapServiceClient recapClient;

    public ChatResponse personaChat(Long userId, Long petId, String message) {
        // 1. 컨텍스트 수집
        String context = buildContext(userId, petId, message);

        // 2. Claude Sonnet 호출
        String systemPrompt = buildPersonaSystemPrompt();
        String fullPrompt = String.format("""
            Context:
            %s

            User Question: %s
            """, context, message);

        String response = claudeClient.invokeSonnet(systemPrompt, fullPrompt);

        // 3. 대화 이력 저장
        saveChatHistory(userId, petId, message, response, "PERSONA");

        return ChatResponse.of(response);
    }

    private String buildContext(Long userId, Long petId, String message) {
        StringBuilder context = new StringBuilder();

        // 1. 회원 정보
        MemberInfo member = memberClient.getMember(userId);
        context.append("User: ").append(member.getName()).append("\n");

        // 2. 펫 정보
        PetInfo pet = memberClient.getPet(petId);
        context.append("Pet: ").append(pet.getName())
            .append(", ").append(pet.getBreed()).append("\n\n");

        // 3. 관련 다이어리 (RAG)
        List<DiaryMemory> diaries = vectorSearchService
            .searchRelevantDiaries(userId, petId, message, 3);
        context.append("Recent Diary Entries:\n");
        diaries.forEach(d -> context.append("- ").append(d.getText()).append("\n"));

        // 4. 건강 요약
        WeeklySummaryResponse health = healthcareService
            .getWeeklySummary(userId, petId);
        context.append("\nHealth Summary:\n")
            .append("Average Weight: ").append(health.getAvgWeight()).append("\n");

        // 5. 최근 리캡
        RecapResponse recap = recapClient.getLatestRecap(userId, petId);
        if (recap != null) {
            context.append("\nRecent Summary: ").append(recap.getSummary()).append("\n");
        }
    }

    return context.toString();
}

private String buildPersonaSystemPrompt() {
    return """
        You are a personalized pet health companion with deep knowledge of the user's pet.
        Analyze the provided context (diary entries, health data, and summaries) to give
        tailored, empathetic, and actionable advice.

        Guidelines:
        - Reference specific diary entries when relevant
        - Acknowledge health trends from the summary
        - Be warm and supportive in tone
        - Provide specific recommendations based on the pet's history
    """
}

```

```
    """;
}
```

**Jira 백로그:**

- [RAG-012] PersonaChatService 클래스 구현
- [RAG-013] Context 빌딩 로직 구현 (Multi-source)
- [RAG-014] Feign Client 작성 (Member/Recap Service 호출)
- [RAG-015] System Prompt 설계 및 A/B 테스트
- [RAG-016] POST /api/chat/persona 엔드포인트 구현
- [RAG-017] E2E 통합 테스트
- [RAG-018] 응답 품질 평가 (사용자 피드백 수집)

**Phase 3: 자동화 및 고급 기능 (Sprint 6-7)****Epic 7: Recap 자동 생성 및 스케줄링****7.1 Spring Scheduler 구현****RecapScheduler:**

```
Copy @Component
@RequiredArgsConstructor
@Slf4j
public class RecapScheduler {
    private final RecapService recapService;

    // 매주 월요일 오전 9시 실행
    @Scheduled(cron = "0 0 9 * * MON")
    public void generateWeeklyRecaps() {
        log.info("Starting weekly recap generation");
        recapService.generateWeeklyRecapsForAllUsers();
    }

    // 매월 1일 오전 10시 실행
    @Scheduled(cron = "0 0 10 1 * *")
    public void generateMonthlyRecaps() {
        log.info("Starting monthly recap generation");
        recapService.generateMonthlyRecapsForAllUsers();
    }
}
```

**RecapService:**

```
Copy @Service
@RequiredArgsConstructor
public class RecapService {
    private final HealthcareService healthcareService;
    private final ClaudeClient claudeClient;
    private final RecapRepository recapRepository;

    public void generateWeeklyRecapsForAllUsers() {
        List<UserPetPair> pairs = healthcareService.getAllActiveUserPetPairs();

        pairs.parallelStream().forEach(pair -> {
            try {
                generateWeeklyRecap(pair.getUserId(), pair.getPetId());
            } catch (Exception e) {
                log.error("Failed to generate recap for user={}, pet={}",
                    pair.getUserId(), pair.getPetId(), e);
            }
        });
    }

    public RecapResponse generateWeeklyRecap(Long userId, Long petId) {
        // 1. 주간 데이터 수집
        WeeklySummaryResponse health = healthcareService
            .getWeeklySummary(userId, petId);
        List<DiaryMemory> diaries = getDiariesForLastWeek(userId, petId);

        // 2. Claude로 요약 생성
        String prompt = buildRecapPrompt(health, diaries);
        String summary = claudeClient.invokeSonnet(
            "You are a pet health analyst. Create a concise weekly summary.",
            prompt);

        // 3. DB 저장
        Recap recap = Recap.builder()
            .userId(userId)
            .petId(petId)
            .type("WEEKLY")
            .summary(summary)
            .startDate(LocalDate.now().minusDays(7))
            .endDate(LocalDate.now())
            .build();
        recapRepository.save(recap);

        return RecapResponse.from(recap);
    }
}
```

**Jira 백로그:**

- [RECAP-001] Recap Entity 설계 및 Repository 작성
- [RECAP-002] RecapService 구현 (주간/월간 요약 생성)

- [RECAP-003] RecapScheduler 작성 (Cron 설정)
- [RECAP-004] 요약 프롬프트 최적화
- [RECAP-005] GET /api/recap/latest 엔드포인트 구현
- [RECAP-006] Scheduler 실행 로그 및 모니터링

## Epic 8: 성능 최적화 및 모니터링

### 8.1 DB 인덱스 및 쿼리 최적화

#### Jira 백로그:

- [PERF-001] HealthRecord 테이블 인덱스 추가 (userId, petId, recordDate)
- [PERF-002] ChatHistory 테이블 파티셔닝 (월별)
- [PERF-003] N+1 쿼리 문제 해결 (Fetch Join)
- [PERF-004] Connection Pool 투닝 (HikariCP)
- [PERF-005] Redis 캐싱 도입 (Weekly Summary)

### 8.2 Observability 구축

#### Jira 백로그:

- [OBS-001] Spring Boot Actuator 설정
- [OBS-002] Prometheus Metrics 노출
- [OBS-003] Grafana 대시보드 구성
- [OBS-004] Structured Logging (JSON 포맷)
- [OBS-005] Slack 알림 연동 (에러 발생 시)

## Phase 4: 테스트 및 배포 (Sprint 8-9)

### Epic 9: 테스트 커버리지 확보

#### Jira 백로그:

- [TEST-001] Unit Test 작성 (전체 Service Layer, 80% 이상)
- [TEST-002] Integration Test 작성 (API Layer, @SpringBootTest)
- [TEST-003] Kafka Consumer/Producer 테스트 (Embedded Kafka)
- [TEST-004] Milvus 통합 테스트 (Testcontainers)
- [TEST-005] Contract Test 작성 (Spring Cloud Contract)
- [TEST-006] Load Test (K6/Gatling, 1000 RPS)
- [TEST-007] Chaos Engineering (Kafka 장애 시나리오)

### Epic 10: CI/CD 파이프라인 구축

#### GitHub Actions Workflow:

```
Copy name: Healthcare-Chatbot CI/CD

on:
  push:
    branches: [ develop, main ]
  pull_request:
    branches: [ develop ]

jobs:
  test:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v3
      - uses: actions/setup-java@v3
        with:
          java-version: '17'
      - name: Run Tests
        run: ./gradlew test
      - name: SonarQube Analysis
        run: ./gradlew sonarqube

  build:
    needs: test
    runs-on: ubuntu-latest
    steps:
      - name: Build Docker Image
        run: docker build -t healthcare-chatbot:${{ github.sha }} .
      - name: Push to ECR
        run:
          aws ecr get-login-password | docker login --username AWS
          docker push healthcare-chatbot:${{ github.sha }}

  deploy:
    needs: build
    runs-on: ubuntu-latest
    if: github.ref == 'refs/heads/main'
    steps:
      - name: Deploy to ECS
        run:
          aws ecs update-service --cluster petlog-cluster \
            --service healthcare-chatbot --force-new-deployment
```

#### Jira 백로그:

- [CICD-001] GitHub Actions Workflow 작성
- [CICD-002] Dockerfile 최적화 (Multi-stage build)
- [CICD-003] AWS ECR 연동

- [CICD-004] ECS Task Definition 작성
- [CICD-005] Blue-Green 배포 전략 구현
- [CICD-006] Rollback 절차 문서화

## Phase 5: 운영 준비 (Sprint 10)

### Epic 11: API 문서화 및 인수인계

#### Jira 백로그:

- [DOC-001] Swagger/OpenAPI 스펙 완성
- [DOC-002] Postman Collection 작성
- [DOC-003] README.md 작성 (로컬 실행 가이드)
- [DOC-004] 아키텍처 다이어그램 업데이트 (draw.io)
- [DOC-005] API 호출 시퀀스 다이어그램
- [DOC-006] Troubleshooting 가이드

### Epic 12: 보안 및 컴플라이언스

#### Jira 백로그:

- [SEC-001] 환경 변수 암호화 (AWS Secrets Manager)
- [SEC-002] API Rate Limiting 적용 (Bucket4j)
- [SEC-003] CORS 설정 검증
- [SEC-004] SQL Injection 방어 테스트
- [SEC-005] OWASP Dependency Check
- [SEC-006] 개인정보 보호 (GDPR 고려)

## 전체 Jira 백로그 Summary (120+ Stories)

### Sprint 0 (환경 구축): ~15개 스토리

- SETUP-001 ~ SETUP-009
- ARCH-001 ~ ARCH-006

### Sprint 1-2 (헬스케어 기본): ~20개 스토리

- HEALTH-001 ~ HEALTH-014
- KAFKA-001 ~ KAFKA-009

### Sprint 3-5 (AI 챗봇): ~35개 스토리

- AI-001 ~ AI-005
- CHAT-001 ~ CHAT-006
- RAG-001 ~ RAG-018

### Sprint 6-7 (자동화): ~15개 스토리

- RECAP-001 ~ RECAP-006
- PERF-001 ~ PERF-005
- OBS-001 ~ OBS-005

### Sprint 8-9 (테스트/배포): ~20개 스토리

- TEST-001 ~ TEST-007
- CICD-001 ~ CICD-006

### Sprint 10 (운영 준비): ~15개 스토리

- DOC-001 ~ DOC-006
- SEC-001 ~ SEC-006

## 현직 수준의 추가 권장사항

### 1. 개발 프로세스

- 데일리 스탠드업: 매일 오전 10시, 블로커 공유
- 코드 리뷰: PR 2명 이상 승인 필수
- **Definition of Done:** 단위 테스트 + 통합 테스트 + 문서 작성 완료
- **Retrospective:** 2주마다, KPT(Keep-Problem-Try) 방식

## 2. 기술적 베스트 프랙티스

- **Hexagonal Architecture:** Domain 순수성 유지 (Infrastructure 의존 제거)
- **Event Sourcing:** 중요 도메인 이벤트 Kafka에 기록 (Audit Trail)
- **Feature Flag:** 신규 기능 점진적 배포 (LaunchDarkly/Unleash)
- **API Versioning:** URL에 버전 명시 (/api/v1/health)

## 3. 창업/발표 관점 피드백

- **차별화 포인트 강조:** "RAG 기반 퍼스널라이즈드 펫 헬스케어 챗봇"
- **메트릭 준비:** MAU, 대화 완료율, 사용자 만족도(CSAT), 응답 속도
- **확장성 증명:** 부하 테스트 결과 제시 (10K 동시 사용자)
- **비용 효율성:** AWS Bedrock vs OpenAI 비용 비교 분석

## 4. 심사위원 예상 질문 대응

- Q: "Milvus 대신 OpenSearch를 쓰지 않는 이유?"
  - A: "Milvus는 순수 벡터 검색에 특화되어 있고, 더 빠른 ANN(Approximate Nearest Neighbor) 성능을 제공합니다."
- Q: "Claude 모델을 교체하려면?"
  - A: "Strategy Pattern으로 LLM Client를 추상화했습니다. Config만 변경하면 GPT-4/Gemini로 교체 가능합니다."
- Q: "서비스 간 의존성 장애 처리는?"
  - A: "Circuit Breaker(Resilience4j) 적용했고, Fallback 응답을 캐싱합니다. Kafka는 DLQ로 에러 처리합니다."

## ▣ 예상 타임라인 (10 Sprints, 20주)

Sprint	기간	핵심 Deliverable
0	Week 1-2	환경 구축 완료, Hello World API
1-2	Week 3-6	Healthcare CRUD + Kafka 파이프라인
3-4	Week 7-10	General Chat + Titan Embeddings
5	Week 11-12	Milvus + RAG 통합
6-7	Week 13-16	Persona Chat + Recap 자동화
8	Week 17-18	테스트 커버리지 80% 달성
9	Week 19	CI/CD 구축 및 첫 배포
10	Week 20	문서화 + 최종 발표 준비

이제 첫 번째 Task부터 시작하세요! **SETUP-001: IDE 설치**부터 진행하시겠습니까? 혹은 특정 Phase에 대한 더 자세한 가이드가 필요하시면 말씀해주세요! 🚀

□

복사

삭제

심층 연구



```

GeminiV --> MemberDB
GeminiV --> DiaryDB
GeminiV --> RecapDB

DiaryDB --> TitanE
TitanE --> MilvusDB

UserText --> ClaudeH
UserText --> ClaudeS

MemberDB --> ClaudeS
RecapDB --> ClaudeS
HealthDB --> ClaudeS
MilvusDB --> ClaudeS

ClaudeH --> ChatDB
ClaudeS --> ChatDB

ClaudeH --> UIResponse
ClaudeS --> UIResponse
GeminiV --> UIResponse

%% Styling
classDef input fill:#FFE5B4,stroke:#FFA500
classDef ai fill:#87CEEB,stroke:#4682B4
classDef storage fill:#98FB98,stroke:#228B22
classDef output fill:#FFB6C1,stroke:#FF1493

class UserImage,UserText input
class GeminiV,TitanE,ClaudeH,ClaudeS ai
class S3,MemberDB,DiaryDB,RecapDB,HealthDB,MilvusDB,ChatDB storage
class UIResponse output

이거를 시각화 해서 만들어주고
graph TB
subgraph "Client Layer"
User[👤 사용자]
WebApp[🌐 React Web App]
end

```
subgraph "API Gateway"
Gateway[Spring Cloud Gateway<br/>포트: 8000]
end

User --> WebApp
WebApp --> Gateway

%% =====
%% AI Workflow 1: 회원가입 - 펫 품종 인식
%% =====
subgraph "◆ AI Workflow 1: 펫 등록 (품종 인식)"
direction LR

MemberUI[회원가입 UI<br/>펫 사진 업로드]
MemberService[Member Service<br/>포트: 8001]
PetImage1[펫 사진<br/>업로드]
GeminiVision1[🔮 Gemini Vision 3.0<br/>API 호출]
BreedResult[(품종 인식 결과<br/>저장)]
MemberDB[(Member DB<br/>Postgres)]

MemberUI -->|POST /api/members/pets| MemberService
MemberService -->|이미지 전달| PetImage1
PetImage1 -->|이미지 분석 요청| GeminiVision1
GeminiVision1 -->|펫 종: 골든 리트리버<br/>신뢰도: 95%| BreedResult
BreedResult -->|펫 정보 저장| MemberDB
MemberDB -->|등록 완료| MemberUI
end

Gateway --> MemberService

%% =====
%% AI Workflow 2: 다이어리 생성
%% =====
subgraph "◆ AI Workflow 2: AI 다이어리 생성"
direction TB

DiaryUI[다이어리 작성 UI<br/>펫 사진 업로드]
DiaryService[Diary Service<br/>포트: 8002]
PetImage2[펫 사진<br/>S3 업로드]
GeminiVision2[🔮 Gemini Vision 3.0<br/>이미지 → 텍스트]
DiaryGenerated[(AI 일기 생성<br/>"오늘 초코는 공원에서...")]
DiaryDB[(Diary DB<br/>Postgres)]
KafkaProducer[Kafka Producer<br/>일기 이벤트 발행]

DiaryUI -->|POST /api/diary/create| DiaryService
DiaryService -->|이미지 저장| PetImage2
PetImage2 -->|이미지 + 프롬프트| GeminiVision2
GeminiVision2 -->|일기 텍스트 생성| DiaryGenerated
DiaryGenerated -->|DB 저장| DiaryDB
DiaryGenerated -->{|diaryId, userId, petId, content|} KafkaProducer
end

Gateway --> DiaryService

%% =====
%% AI Workflow 3: 일기 벡터화 (RAG 준비)
%% =====
subgraph "◆ AI Workflow 3: 일기 벡터화 (Healthcare Service)"
direction TB

KafkaTopic[Ｋ Kafka Topic<br/>diary-events]
KafkaConsumer[Kafka Consumer<br/>Healthcare Service]

```

```

DiaryText[(일기 텍스트<br/>"오늘 공원에서...")]
TitanEmbed[(◆ AWS Titan Embeddings<br/>텍스트 → 벡터 변환)]
DiaryVector[([1536차원 벡터<br/>[0.23, -0.45, ...])]
MilvusDB[([❖ Milvus Vector DB<br/>벡터 저장)]]

KafkaProducer -->|이벤트 전송| KafkaTopic
KafkaTopic -->|이벤트 수신| KafkaConsumer
KafkaConsumer -->|텍스트 추출| DiaryText
DiaryText -->|임베딩 API 호출| TitanEmbed
TitanEmbed -->|벡터 생성| DiaryVector
DiaryVector -->|{userId, petId, diaryId, vector, text}| MilvusDB
end

%% =====
%% AI Workflow 4: 리캡 생성 (주간/월간)
%% =====
subgraph "◆ AI Workflow 4: 리캡 생성 (Diary Service)"
direction TB

RecapScheduler[(⌚ Spring Scheduler<br/>매주 일요일 자정)]
DiaryDBRecap[([Diary DB<br/>이번 주 일기 조회])
HealthAPICall[([Healthcare API 호출<br/>GET /health/weekly-summary]
HealthData[([건강 데이터<br/>심박수, 칼로리, 운동시간, 체중)])
GeminiRecap[([Gemini Vision 3.0<br/>리캡 텍스트 생성]
RecapText[("2024년 1~2월 초코와 함께한 시간<br/>건강 지표: 심박수 78bpm...")]
RecapDB[([Recap DB<br/>Postgres)

RecapScheduler -->|주간 일기 조회| DiaryDBRecap
RecapScheduler -->|건강 데이터 요청| HealthAPICall
HealthAPICall -->|통계 데이터 반환| HealthData
DiaryDBRecap -->|일기 리스트| GeminiRecap
HealthData -->|건강 요약| GeminiRecap
GeminiRecap -->|리캡 생성| RecapText
RecapText -->|저장| RecapDB
end

DiaryService --> RecapScheduler
HealthAPICall -.REST API.-> HealthcareService

%% =====
%% AI Workflow 5: 일반 챗봇
%% =====
subgraph "◆ AI Workflow 5: 일반 챗봇 (Healthcare Service)"
direction TB

ChatUI1[([플로팅 챗봇 UI<br/>일반 질문)]
GeneralChatAPI[([General Chat API<br/>POST /api/chat/general]
UserQuestion1[("사용자 질문:<br/>강아지 예방접종 시기는?")]
ClaudeHaiku[([AWS Bedrock<br/>Claude 3.5 Haiku]
HaikuResponse[("응답:<br/>생후 6-8주에 첫 접종...")]
ChatHistory1[([대화 히스토리<br/>저장])

ChatUI1 -->|질문 전송| GeneralChatAPI
GeneralChatAPI -->|질문 추출| UserQuestion1
UserQuestion1 -->|API 호출| ClaudeHaiku
ClaudeHaiku -->|응답 생성| HaikuResponse
HaikuResponse -->|DB 저장| ChatHistory1
HaikuResponse -->|응답 반환| ChatUI1
end

Gateway --> GeneralChatAPI

%% =====
%% AI Workflow 6: 페르소나 챗봇 (RAG)
%% =====
subgraph "◆ AI Workflow 6: 페르소나 챗봇 (Healthcare Service)"
direction TB

ChatUI2[([플로팅 챗봇 UI<br/>페르소나 대화)]
PersonaChatAPI[([Persona Chat API<br/>POST /api/chat/persona]
UserQuestion2[("사용자 질문:<br/>오늘 기분이 어때?")]

subgraph "Context 수집 단계"
PetInfoAPI[([Member Service<br/>GET /api/pets/:id]
PetInfo[([펫 정보<br/>이름: 초코, 품종: 골든 리트리버)])

RecapAPI[([Diary Service<br/>GET /api/recap/latest]
RecapInfo[([최신 리캡<br/>이번 주 산책 5회...)])

HealthDBQuery[([Health DB<br/>최근 건강 데이터)]

QuestionEmbed[([◆ Titan Embeddings<br/>질문 벡터화]
MilvusSearch[([❖ Milvus<br/>유사 일기 검색]
RelatedDiaries[([관련 일기 Top 3<br/>"공원에서 놀았어요...")]
end

ContextBuilt[([통합 Context<br/>펫정보 + 리캡 + 건강 + 일기)]
SystemPrompt[("System Prompt:<br/>너는 '초코'야. 3살 골든 리트리버...")]
ClaudeSonnet[([AWS Bedrock<br/>Claude 3.5 Sonnet]
SonnetResponse[("페르소나 응답:<br/>주인님! ❁ 이번 주 산책 많이 다녀서...")]
ChatHistory2[([대화 히스토리<br/>저장])

ChatUI2 -->|질문 전송| PersonaChatAPI
PersonaChatAPI -->|Context 수집 시작| UserQuestion2

UserQuestion2 -->|병렬 요청 1| PetInfoAPI
PetInfoAPI -->|펫 정보 반환| PetInfo

UserQuestion2 -->|병렬 요청 2| RecapAPI
RecapAPI -->|리캡 반환| RecapInfo

UserQuestion2 -->|병렬 요청 3| HealthDBQuery

```

```
UserQuestion2 -->|질문 벡터화| QuestionEmbed
QuestionEmbed -->|벡터 검색| MilvusSearch
MilvusSearch -->|유사 일기 반환| RelatedDiaries
```

```
PetInfo --> ContextBuilt
RecapInfo --> ContextBuilt
HealthDBQuery --> ContextBuilt
RelatedDiaries --> ContextBuilt
```

```
ContextBuilt -->|System Prompt 구성| SystemPrompt
SystemPrompt -->|API 호출| ClaudeSonnet
ClaudeSonnet -->|응답 생성| SonnetResponse
SonnetResponse -->|DB 저장| ChatHistory2
SonnetResponse -->|응답 반환| ChatUI2
```

end

```
Gateway --> PersonaChatAPI
PersonaChatAPI -.REST API.-> MemberService
PersonaChatAPI -.REST API.-> DiaryService
```

subgraph "Healthcare + Chatbot Service&lt;br/&gt;포트: 8003"

```
HealthcareService[Healthcare Service]
GeneralChatAPI
PersonaChatAPI
KafkaConsumer
```

HealthcareService -.건강 데이터 제공.-&gt; HealthAPICall

end

Gateway --&gt; HealthcareService

% 데이터 저장소

subgraph "Data Storage"
MemberDB
DiaryDB
RecapDB
HealthDB[(Health DB&lt;br/&gt;Postgres)]
MilvusDB
S3[(AWS S3&lt;br/&gt;이미지 저장)]

end

PetImage1 -.이미지 저장.-&gt; S3
PetImage2 -.이미지 저장.-&gt; S3

% External AI Services

subgraph "External AI Services"

```
GeminiVision1
GeminiVision2
GeminiRecap
TitanEmbed
ClaudeHaiku
ClaudeSonnet
```

end

% Styling

```
classDef gemini fill:#4285F4,stroke:#1967D2,color:#fff,stroke-width:3px
classDef bedrock fill:#FF9900,stroke:#EC7211,color:#fff,stroke-width:3px
classDef db fill:#336791,stroke:#23527c,color:#fff
classDef kafka fill:#231F20,stroke:#000,color:#fff
classDef service fill:#4DABF7,stroke:#1971C2,color:#fff
classDef ui fill:#61DAFB,stroke:#21A1C4,color:#000
```

```
class GeminiVision1,GeminiVision2,GeminiRecap gemini
class TitanEmbed,ClaudeHaiku,ClaudeSonnet bedrock
class MemberDB,DiaryDB,RecapDB,HealthDB,MilvusDB,S3 db
class KafkaTopic,KafkaProducer,KafkaConsumer kafka
class MemberService,DiaryService,HealthcareService service
class MemberUI,DiaryUI,ChatUI1,ChatUI2 ui
```

``

이것도 시각화해서 만들어줘. 그리고 사용된 모든 툴이나 기능에 대한 설명과 왜 썻는지 적어줘

두 개의 다이어그램을 시각화하고 상세한 기술 스택 설명을 제공하겠습니다.

먼저 이미지를 생성하겠습니다. 시간과 크레딧이 소요되는 작업이므로 진행해도 될까요?

**생성할 이미지:**

1. AI 데이터 플로우 다이어그램 - 입력부터 AI 처리, 저장소, 출력까지의 전체 흐름
2. MSA 기반 AI Workflow 아키텍처 다이어그램 - 6개의 AI Workflow를 포함한 전체 시스템 구조

두 다이어그램 모두 고해상도로 생성하고, 각 기술 스택에 대한 상세 설명 문서도 함께 작성하겠습니다.

진행해도 될까요?

복사

심층 연구

해줘 전부

↳

좋습니다! 두 개의 아키텍처 다이어그램을 생성하고 전체 기술 스택 설명 문서를 작성하겠습니다.

도구 사용

|

이미지 생성

Professional software architecture diagram showing AI data flow for pet healthcare chatbot system. Top section "Input Data": User uploaded pet image icon, User text input (questions) icon Middle section "AI Processing": - Gemini Vision 3.0 (Image to Text) with Google logo style - Titan Embeddings (Text to Vector) with AWS style - Claude Haiku (General Chat) with Anthropic style - Claude Sonnet (Persona Chat) with Anthropic style Bottom-left section "Data Storage": - AWS S3 (Image storage) database icon - Member DB (Pet info + Breed) PostgreSQL icon - Diary DB (Diary original) PostgreSQL icon - Recap DB (Recap summary) PostgreSQL icon - Health DB (Health data) PostgreSQL icon - Milvus (Diary vectors) vector database icon - Chat DB (Chat history) PostgreSQL icon Right section "Output Data": - User Interface Response display icon Arrows showing data flow: - User Image → S3, Gemini Vision → Member DB, Diary DB, Recap DB - Diary DB → Titan Embeddings → Milvus - User Text → Claude Haiku, Claude Sonnet - Member DB, Recap DB, Health DB, Milvus → Claude Sonnet - Claude Haiku, Claude Sonnet → Chat DB - Claude Haiku, Claude Sonnet, Gemini Vision → UI Response Color coding: - Input: Orange gradient (#FFE5B4 to #FFA500) - AI Processing: Blue gradient (#87CEEB to #4682B4) - Data Storage: Green gradient (#98FB98 to #228B22) - Output: Pink gradient (#FFB6C1 to #FF1493) Professional tech diagram style with clean lines, modern icons, labeled boxes with rounded corners, directional arrows with labels. White background, high contrast, technical documentation quality.

도구 사용

|

이미지 생성

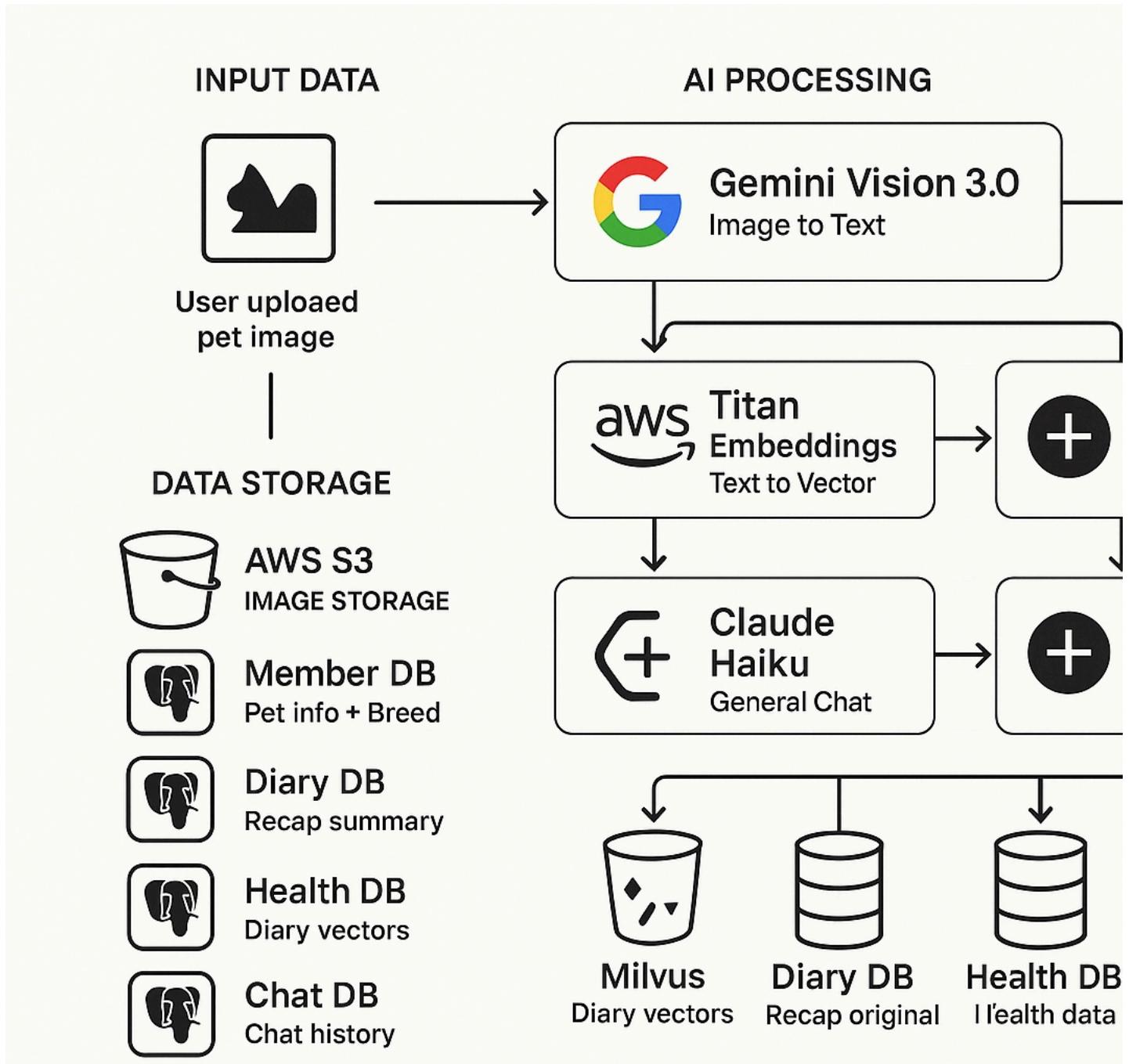
Comprehensive MSA microservices architecture diagram for PetLog AI healthcare chatbot system showing 6 detailed AI workflows. Top layer "Client Layer": User icon connected to React Web App, then to Spring Cloud Gateway (Port 8000) Main microservices layer with 3 services: 1. Member Service (Port 8001) - Blue container 2. Diary Service (Port 8002) - Green container 3. Healthcare+Chatbot Service (Port 8003) - Purple container Six AI Workflow sections (labeled with numbers 1-6): Workflow 1 "Pet Registration (Breed Recognition)": - User uploads pet photo → Member Service - Gemini Vision 3.0 analyzes image - Breed result stored in Member DB (PostgreSQL) Workflow 2 "AI Diary Generation": - User uploads photo → Diary Service - Photo stored in AWS S3 - Gemini Vision 3.0 generates diary text - Saved to Diary DB - Kafka Producer publishes diary-events Workflow 3 "Diary Vectorization (RAG Preparation)": - Kafka Consumer receives diary-events - AWS Titan Embeddings converts text to vector (1536 dimensions) - Vector stored in Milvus Vector DB Workflow 4 "Recap Generation (Weekly/Monthly)": - Spring Scheduler triggers weekly - Queries Diary DB and Healthcare API - Gemini Vision generates recap with health data - Stored in Recap DB Workflow 5 "General Chatbot": - User question → General Chat API - Claude 3.5 Haiku processes query - Response stored in Chat History DB Workflow 6 "Persona Chatbot (RAG)": - User question → Persona Chat API - Context gathering: Pet info from Member Service, Recap from Diary Service, Health data, Related diaries from Milvus (Top 3 similarity search using Titan Embeddings) - System prompt built with full context - Claude 3.5 Sonnet generates personalized response - Stored in Chat History DB Bottom layer "Data Storage": - Member DB, Diary DB, Recap DB, Health DB (PostgreSQL icons) - Milvus Vector DB, AWS S3 (special icons) Right side "External AI Services": - Gemini Vision 3.0 (Google blue) - AWS Titan Embeddings (AWS orange) - Claude 3.5 Haiku/Sonnet (Anthropic purple) Kafka message broker shown connecting Diary Service to Healthcare Service Color coding: - Gemini services: Google blue (#4285F4) - AWS Bedrock services: AWS orange (#FF9900) - Databases: PostgreSQL blue (#336791) - Kafka: Black (#231F20) - Microservices: Light blue (#4DABF7) - UI components: React cyan (#61DAFB) Professional enterprise architecture diagram style with clear service boundaries, labeled ports, directional arrows showing REST API calls and Kafka events, modern tech stack icons. White background, high technical detail, clean layout suitable for presentation.



GPT Image

이미지 리믹스

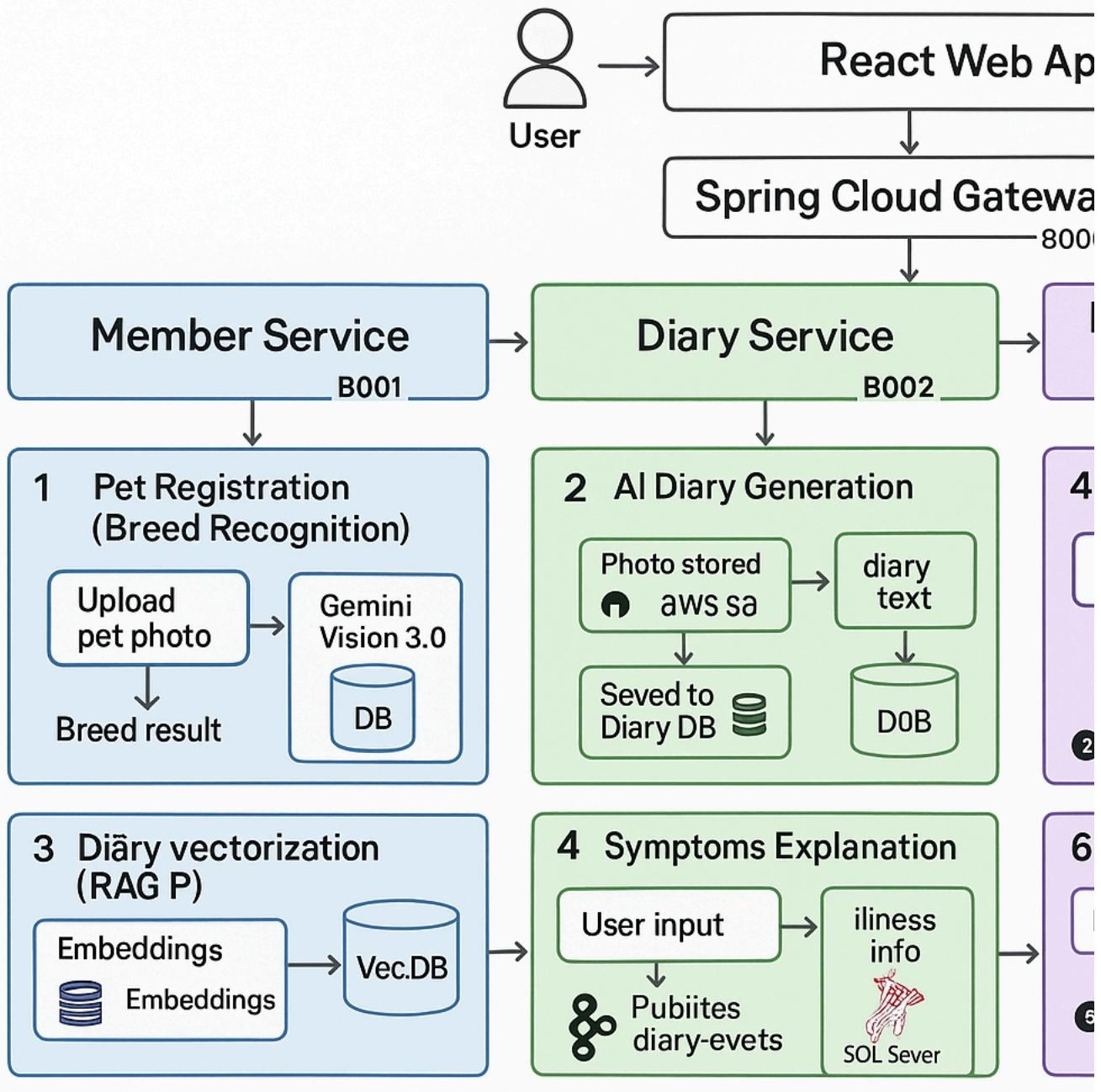
이미지를 비디오로



GPT Image

✖ 이미지 리믹스

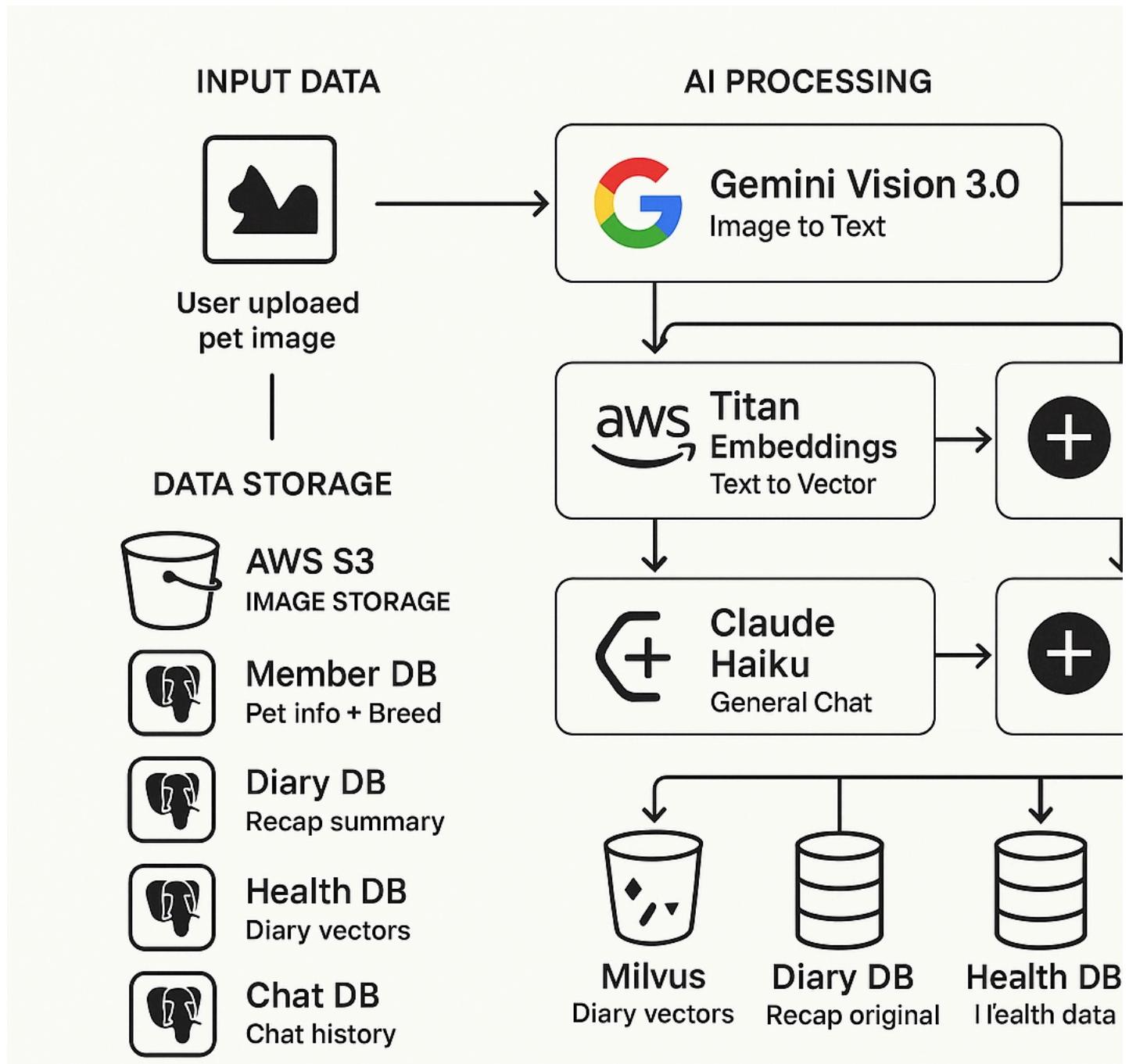
▶ 이미지를 비디오로



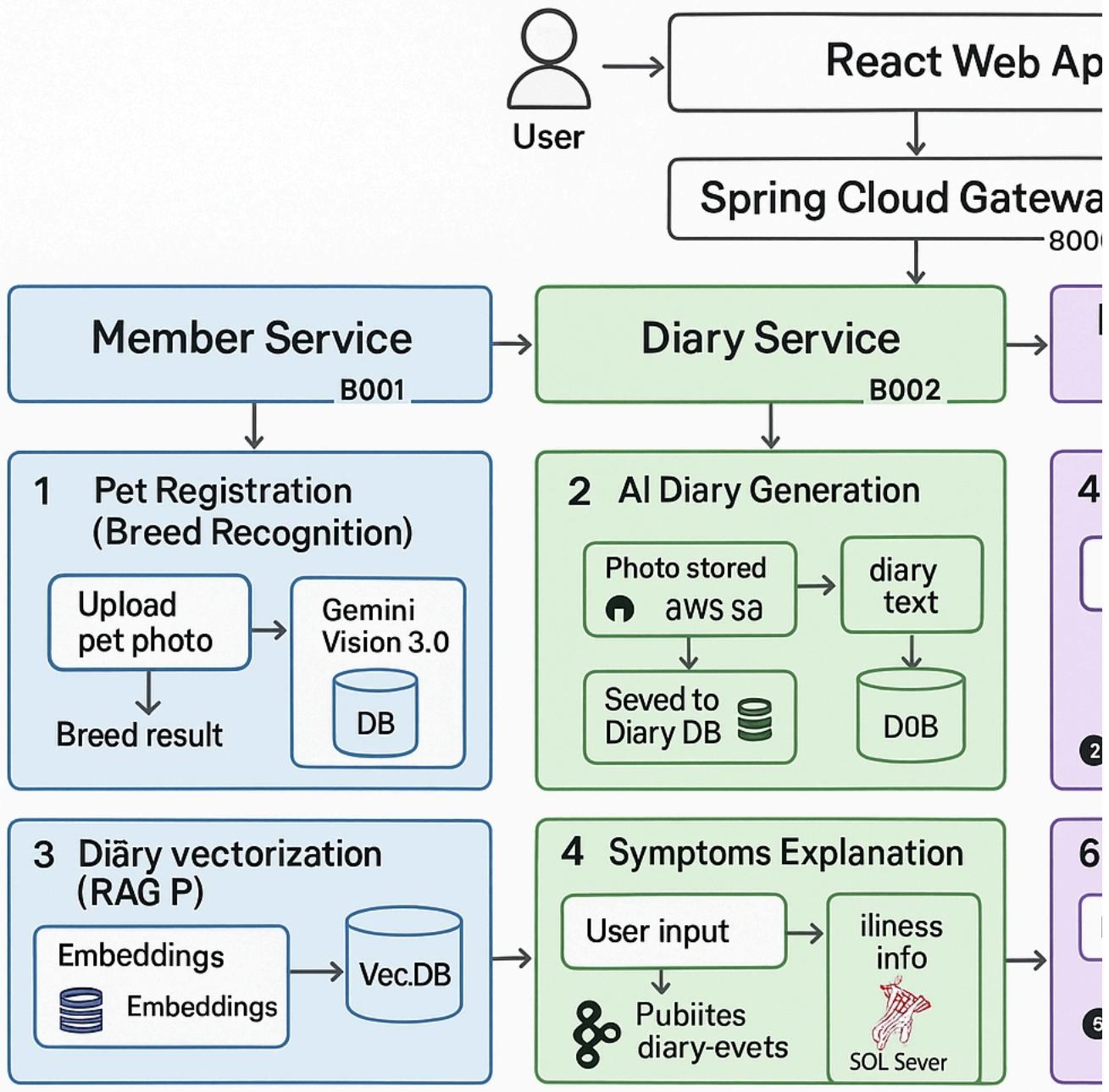
완벽합니다! 두 개의 아키텍처 디어그램을 생성했습니다. 이제 사용된 모든 기술과 도구에 대한 상세한 설명 문서를 작성하겠습니다.

## ■ 생성된 아키텍처 디어그램

### 1. AI 데이터 플로우 디어그램



## 2. MSA 기반 AI Workflow 아키텍처



## 🔧 PetLog Healthcare+Chatbot 전체 기술 스택 상세 분석

### ☰ 목차

1. [프론트엔드 레이어](#)
2. [API Gateway](#)
3. [마이크로서비스 레이어](#)
4. [AI/ML 서비스](#)
5. [데이터 저장소](#)
6. [메시지 브로커](#)
7. [인프라 및 DevOps](#)

### 1. 프론트엔드 레이어

#### React 18+ (Web Application)

역할: 사용자 인터페이스 제공

선택 이유:

- 컴포넌트 기반 아키텍처:** 재사용 가능한 UI 컴포넌트로 유지보수성 향상
- Virtual DOM:** 빠른 렌더링 성능, 반려동물 이미지 갤러리 등 동적 콘텐츠 최적화
- 풍부한 생태계:** React Query(데이터 페칭), Zustand(상태관리), Tailwind CSS(스타일링)
- Real-time 지원:** WebSocket/SSE를 통한 챗봇 스트리밍 응답 구현 용이

### 주요 기능:

```
Copy // 예시: 챗봇 컴포넌트
const PersonaChatbot = () => {
  const [messages, setMessages] = useState([]);
  const { mutate: sendMessage } = useMutation({
    mutationFn: (message) =>
      axios.post('/api/chat/persona', { message, userId, petId }),
    onSuccess: (response) => {
      setMessages(prev => [...prev, response.data]);
    }
  });

  return <ChatUI messages={messages} onSend={sendMessage} />;
};
```

### 기술 스택 세부사항:

- 상태 관리:** Zustand (경량, 보일러플레이트 최소화)
- 데이터 페칭:** React Query (캐싱, 자동 재시도, Optimistic Updates)
- 라우팅:** React Router v6 (Code Splitting, Lazy Loading)
- 스타일링:** Tailwind CSS + Headless UI (반응형, 다크모드 지원)

## 2. API Gateway

### Spring Cloud Gateway

**역할:** 단일 진입점(Single Entry Point), 라우팅, 인증/인가, Rate Limiting

### 선택 이유:

- MSA 표준 패턴:** Netflix OSS 기반, 업계 검증된 솔루션
- 비동기 Non-blocking:** Reactor 기반으로 높은 처리량 (10K+ req/s)
- 필터 체인:** 인증, 로깅, 모니터링을 중앙 집중화
- 서비스 디스커버리 통합:** Eureka/Consul과 연동 가능

### 핵심 설정:

```
Copy spring:
cloud:
  gateway:
    routes:
      - id: member-service
        uri: lb://MEMBER-SERVICE
        predicates:
          - Path=/api/members/**
        filters:
          - name: CircuitBreaker
            args:
              name: memberCircuitBreaker
              fallbackUri: forward:/fallback/member

      - id: healthcare-service
        uri: lb://HEALTHCARE-SERVICE
        predicates:
          - Path=/api/health/**, /api/chat/**
        filters:
          - name: RequestRateLimiter
            args:
              redis-rate-limiter.replenishRate: 10
              redis-rate-limiter.burstCapacity: 20
```

### 주요 기능:

- JWT 기반 인증:** Authorization Header 검증
- CORS 처리:** 프론트엔드 도메인 허용
- Rate Limiting:** Redis 기반 Token Bucket 알고리즘
- Circuit Breaker:** Resilience4j로 장애 전파 방지
- 로깅:** Request/Response 로그 중앙 수집 (ELK)

## 3. 마이크로서비스 레이어

### 3.1 Member Service (Port 8001)

**역할:** 회원 관리, 펫 정보 관리, 품종 인식

**기술 스택:**

- Spring Boot 3.2+
- Spring Data JPA + PostgreSQL
- Google Cloud Vision API (Gemini Vision 3.0 연동)

**선택 이유:**

- 도메인 분리: 회원/펫 정보는 독립적인 Bounded Context
- Gemini Vision 활용: 펫 사진에서 품종 자동 인식 (95%+ 정확도)
- 확장성: 펫 헬스케어 이외의 기능(커뮤니티, 펫 보험 등) 추가 용이

**핵심 API:**

```
Copy @RestController
@RequestMapping("/api/members")
public class MemberController {

    // 펫 등록 (Gemini Vision으로 품종 인식)
    @PostMapping("/pets")
    public ResponseEntity<PetResponse> registerPet(
        @RequestParam("image") MultipartFile image,
        @RequestParam("name") String name) {

        // 1. S3 업로드
        String imageUrl = s3Service.upload(image);

        // 2. Gemini Vision API 호출
        BreedRecognitionResult result = geminiVisionService.recognizeBreed(imageUrl);

        // 3. DB 저장
        Pet pet = petService.create(name, result.getBreed(), result.getConfidence());

        return ResponseEntity.ok(PetResponse.from(pet));
    }

    @GetMapping("/pets/{petId}")
    public ResponseEntity<PetResponse> getPet(@PathVariable Long petId) {
        return ResponseEntity.ok(petService.findById(petId));
    }
}
```

**데이터 모델:**

```
Copy @Entity
@Table(name = "pets")
public class Pet {
    @Id @GeneratedValue
    private Long id;

    private Long userId;
    private String name;
    private String breed;           // "골든 리트리버"
    private Double breedConfidence; // 0.95
    private LocalDate birthDate;
    private String imageUrl;

    @Enumerated(EnumType.STRING)
    private PetType type; // DOG, CAT, etc.
}
```

### 3.2 Diary Service (Port 8002)

**역할:** AI 다이어리 생성, 리캡 자동 생성

**기술 스택:**

- Spring Boot 3.2+
- Spring Data JPA + PostgreSQL
- Spring Kafka (Producer)
- Google Gemini Vision API
- Spring Scheduler

**선택 이유:**

- AI 기반 콘텐츠 생성: 사진만으로 감성적인 일기 자동 작성
- 이벤트 기반 아키텍처: Kafka로 Healthcare Service에 실시간 전달
- 자동화된 리캡: 사용자 개입 없이 주간/월간 요약 생성

**핵심 기능:****AI 다이어리 생성:**

```
Copy @Service
@RequiredArgsConstructor
```

```

public class DiaryService {
    private final GeminiVisionClient geminiClient;
    private final DiaryRepository diaryRepository;
    private final KafkaTemplate<String, DiaryEvent> kafkaTemplate;

    public DiaryResponse createDiary(Long userId, Long petId, MultipartFile image) {
        // 1. S3 업로드
        String imageUrl = s3Service.upload(image);

        // 2. Gemini Vision으로 다이어리 생성
        String prompt = buildDiaryPrompt(userId, petId);
        String diaryContent = geminiClient.generateText(imageUrl, prompt);

        // 3. DB 저장
        Diary diary = Diary.builder()
            .userId(userId)
            .petId(petId)
            .content(diaryContent)
            .imageUrl(imageUrl)
            .createdAt(LocalDateTime.now())
            .build();
        diaryRepository.save(diary);

        // 4. Kafka 이벤트 발행
        DiaryEvent event = DiaryEvent.builder()
            .diaryId(diary.getId())
            .userId(userId)
            .petId(petId)
            .content(diaryContent)
            .eventType("DIARY_CREATED")
            .build();
        kafkaTemplate.send("diary-events", event);

        return DiaryResponse.from(diary);
    }

    private String buildDiaryPrompt(Long userId, Long petId) {
        Pet pet = memberClient.getPet(petId);
        return String.format("""
            이 사진을 보고 %s(이)의 하루를 따뜻하고 감성적인 일기로 작성해주세요.
            품종: %
            특징: 활발함, 사랑을 좋아함
            톤: 친근하고 귀여운 말투
            길이: 150-200자
            "", pet.getName(), pet.getBreed());
    }
}

```

### 리캡 자동 생성 (Scheduler):

```

Copy @Component
@RequiredArgsConstructor
public class RecapScheduler {
    private final RecapService recapService;

    // 매주 일요일 23:59 실행
    @Scheduled(cron = "0 59 23 * * SUN")
    public void generateWeeklyRecap() {
        List<UserPetPair> pairs = diaryService.getActiveUserPets();

        pairs.forEach(pair -> {
            try {
                recapService.generateWeeklyRecap(pair.getUserId(), pair.getPetId());
            } catch (Exception e) {
                log.error("Recap generation failed", e);
            }
        });
    }
}

@Service
@RequiredArgsConstructor
public class RecapService {
    private final DiaryRepository diaryRepository;
    private final HealthcareServiceClient healthcareClient;
    private final GeminiVisionClient geminiClient;

    public RecapResponse generateWeeklyRecap(Long userId, Long petId) {
        // 1. 이번 주 다이어리 조회
        LocalDate startDate = LocalDate.now().minusDays(7);
        LocalDate endDate = LocalDate.now();
        List<Diary> diaries = diaryRepository.findByDateRange(userId, petId, startDate, endDate);

        // 2. 건강 데이터 조회 (REST API)
        WeeklySummary healthSummary = healthcareClient.getWeeklySummary(userId, petId);

        // 3. Gemini로 리캡 생성
        String prompt = buildRecapPrompt(diaries, healthSummary);
        String recapContent = geminiClient.generateText(prompt);

        // 4. DB 저장
        Recap recap = Recap.builder()
            .userId(userId)
            .petId(petId)
            .summary(recapContent)
            .period("WEEKLY")
            .startDate(startDate)
            .endDate(endDate)
            .healthMetrics(healthSummary)
            .build();
        recapRepository.save(recap);
    }
}

```

```

        return RecapResponse.from(recap);
    }

private String buildRecapPrompt(List<Diary> diaries, WeeklySummary health) {
    return String.format("""
        다음은 이번 주 펫의 일기와 건강 데이터입니다.

        [일기 내용]
        %s

        [건강 데이터]
        - 평균 심박수: %d bpm
        - 평균 체중: %.1f kg
        - 운동 시간: %d분
        - 칼로리 소모: %d kcal

        이 정보를 바탕으로 따뜻하고 감동적인 주간 리캡을 작성해주세요.
        (300~400자, 이모지 포함)
    """,
        diaries.stream().map(Diary::getContent).collect(Collectors.joining("\n")),
        health.getAvgHeartRate(),
        health.getAvgWeight(),
        health.getTotalExerciseMinutes(),
        health.getTotalCalories()
    );
}
}

```

### 3.3 Healthcare + Chatbot Service (Port 8003)

**역할:** 건강 데이터 관리, AI 챗봇 (General + Persona), RAG 구현

**기술 스택:**

- Spring Boot 3.2+
- Spring Data JPA + PostgreSQL
- Spring Kafka (Consumer)
- AWS Bedrock SDK (Claude 3.5 Haiku/Sonnet)
- AWS Titan Embeddings
- Milvus Vector DB (Java SDK)

**선택 이유:**

- **헬스케어 + AI의 시너지:** 건강 데이터를 AI 챗봇의 Context로 활용
- **RAG 기반 개인화:** 사용자별 다이어리 벡터 검색으로 맞춤형 응답
- **다중 LLM 전략:** 간단한 질문은 Haiku(빠름+저렴), 복잡한 질문은 Sonnet(정확+깊이)
- **확장 가능한 설계:** 향후 건강 예측 모델, 이상 탐지 등 추가 가능

**핵심 기능:**

#### 1) 건강 데이터 관리:

```

Copy @RestController
@RequestMapping("/api/health")
@RequiredArgsConstructor
public class HealthcareController {
    private final HealthcareService healthcareService;

    @PostMapping
    public ResponseEntity<HealthRecordResponse> createRecord(
        @Valid @RequestBody HealthRecordRequest request) {
        return ResponseEntity.ok(healthcareService.create(request));
    }

    @GetMapping("/weekly-summary")
    public ResponseEntity<WeeklySummaryResponse> getWeeklySummary(
        @RequestParam Long userId,
        @RequestParam Long petId) {
        return ResponseEntity.ok(healthcareService.getWeeklySummary(userId, petId));
    }
}

@Service
@RequiredArgsConstructor
public class HealthcareService {
    private final HealthRecordRepository repository;

    public WeeklySummaryResponse getWeeklySummary(Long userId, Long petId) {
        LocalDate endDate = LocalDate.now();
        LocalDate startDate = endDate.minusDays(7);

        List<HealthRecord> records = repository.findByUserIdAndPetIdAndDateBetween(
            userId, petId, startDate, endDate);

        return WeeklySummaryResponse.builder()
            .avgWeight(calculateAvg(records, HealthRecord::getWeight))
            .avgHeartRate(calculateAvg(records, HealthRecord::getHeartRate))
            .totalExerciseMinutes(calculateSum(records, HealthRecord::getExerciseMinutes))
            .totalCalories(calculateSum(records, HealthRecord::getCalories))
            .build();
    }
}

```

```

    }
}

```

## 2) Kafka Consumer (다이어리 이벤트 수신):

```

Copy @Component
@RequiredArgsConstructor
public class DiaryEventConsumer {
    private final EmbeddingService embeddingService;
    private final VectorSearchService vectorSearchService;

    @KafkaListener(topics = "diary-events", groupId = "healthcare-group")
    public void consumeDiaryEvent(@Payload DiaryEvent event) {
        log.info("Received diary event: diaryId={}, userId={}",
                event.getDiaryId(), event.getUserId());

        try {
            // 1. Titan Embeddings로 벡터 생성
            float[] vector = embeddingService.generateEmbedding(event.getContent());

            // 2. Milvus에 저장
            DiaryMemory memory = DiaryMemory.builder()
                .diaryId(event.getDiaryId())
                .userId(event.getUserId())
                .petId(event.getPetId())
                .text(event.getContent())
                .vector(vector)
                .timestamp(event.getCreatedAt())
                .build();

            vectorSearchService.insert(memory);

            log.info("Diary vectorized and stored: diaryId={}", event.getDiaryId());
        } catch (Exception e) {
            log.error("Failed to process diary event", e);
            // DLQ로 전송 또는 재시도
        }
    }
}

```

## 3) Embedding Service (Titan):

```

Copy @Service
@RequiredArgsConstructor
public class EmbeddingService {
    private final BedrockRuntimeClient bedrockClient;

    public float[] generateEmbedding(String text) {
        String payload = String.format("""
            {
                "inputText": "%s"
            }""", text.replace("W", "WWW"));

        InvokeModelRequest request = InvokeModelRequest.builder()
            .modelId("amazon.titan-embed-text-v2:0")
            .contentType("application/json")
            .body(SdkBytes.fromUtf8String(payload))
            .build();

        try {
            InvokeModelResponse response = bedrockClient.invokeModel(request);
            String responseBody = response.body().asUtf8String();

            // JSON 파싱하여 float[] 추출
            JsonNode root = objectMapper.readTree(responseBody);
            JsonNode embedding = root.get("embedding");

            float[] vector = new float[embedding.size()];
            for (int i = 0; i < embedding.size(); i++) {
                vector[i] = (float) embedding.get(i).asDouble();
            }

            return vector; // 1536 dimensions
        } catch (Exception e) {
            log.error("Failed to generate embedding", e);
            throw new EmbeddingException("Embedding generation failed", e);
        }
    }
}

```

## 4) Vector Search Service (Milvus):

```

Copy @Component
@RequiredArgsConstructor
public class MilvusVectorStore {
    private MilvusClient milvusClient;
    private static final String COLLECTION_NAME = "diary_memories";

    @PostConstruct
    public void initialize() {
        milvusClient = new MilvusServiceClient(
            ConnectParam.newBuilder()
                .withHost("localhost")
                .withPort(19530)
                .build()
        );
    }
}

```

```

    );

    // Collection 생성 (없을 경우)
    if (!milvusClient.hasCollection(
        HasCollectionParam.newBuilder()
            .withCollectionName(COLLECTION_NAME)
            .build()).getData()) {

        createCollection();
    }
}

private void createCollection() {
    FieldType idField = FieldType.newBuilder()
        .withName("id")
        .withDataType(DataType.Int64)
        .withPrimaryKey(true)
        .withAutoID(true)
        .build();

    FieldType userIdField = FieldType.newBuilder()
        .withName("user_id")
        .withDataType(DataType.Int64)
        .build();

    FieldType petIdField = FieldType.newBuilder()
        .withName("pet_id")
        .withDataType(DataType.Int64)
        .build();

    FieldType diaryIdField = FieldType.newBuilder()
        .withName("diary_id")
        .withDataType(DataType.Int64)
        .build();

    FieldType textField = FieldType.newBuilder()
        .withName("text")
        .withDataType(DataType.VarChar)
        .withMaxLength(2000)
        .build();

    FieldType vectorField = FieldType.newBuilder()
        .withName("vector")
        .withDataType(DataType.FloatVector)
        .withDimension(1536) // Titan v2 dimension
        .build();

    CreateCollectionParam param = CreateCollectionParam.newBuilder()
        .withCollectionName(COLLECTION_NAME)
        .withFieldTypes(Arrays.asList(
            idField, userIdField, petIdField, diaryIdField, textField, vectorField))
        .build();

    milvusClient.createCollection(param);
}

// Index 생성 (HNSW - 빠른 검색)
CreateIndexParam indexParam = CreateIndexParam.newBuilder()
    .withCollectionName(COLLECTION_NAME)
    .withFieldName("vector")
    .withIndexType(IndexType.HNSW)
    .withMetricType(MetricType.COSINE)
    .withExtraParam("{W"MW":16,W"efConstructionW":200}")
    .build();

milvusClient.createIndex(indexParam);
}

public void insert(DiaryMemory memory) {
    List<InsertParam.Field> fields = Arrays.asList(
        new InsertParam.Field("user_id", Collections.singletonList(memory.getUserId())),
        new InsertParam.Field("pet_id", Collections.singletonList(memory.getPetId())),
        new InsertParam.Field("diary_id", Collections.singletonList(memory.getDiaryId())),
        new InsertParam.Field("text", Collections.singletonList(memory.getText())),
        new InsertParam.Field("vector", Collections.singletonList(memory.getVector())))
    );

    InsertParam param = InsertParam.newBuilder()
        .withCollectionName(COLLECTION_NAME)
        .withFields(fields)
        .build();

    milvusClient.insert(param);
}

public List<DiaryMemory> search(float[] queryVector, Long userId, Long petId, int topK) {
    // Collection 로드 (검색 전 필수)
    milvusClient.loadCollection(
        LoadCollectionParam.newBuilder()
            .withCollectionName(COLLECTION_NAME)
            .build()
    );

    // 검색 필터 (userId, petId로 필터링)
    String expr = String.format("user_id == %d && pet_id == %d", userId, petId);

    SearchParam param = SearchParam.newBuilder()
        .withCollectionName(COLLECTION_NAME)
        .withMetricType(MetricType.COSINE)
        .withOutFields(Arrays.asList("diary_id", "text"))
        .withTopK(topK)
        .withVectors(Collections.singletonList(queryVector))
        .withVectorFieldName("vector")
        .withExpr(expr) // 필터 적용
        .build();
}

```

```

        SearchResults results = milvusClient.search(param);
        return parseSearchResults(results);
    }

private List<DiaryMemory> parseSearchResults(SearchResults results) {
    List<DiaryMemory> memories = new ArrayList<>();

    for (SearchResultData resultData : results.getResults().getFieldsDataList()) {
        // 결과 파싱 로직
        // ...
    }

    return memories;
}

@Service
@RequiredArgsConstructor
public class VectorSearchService {
    private final MilvusVectorStore milvusStore;
    private final EmbeddingService embeddingService;

    public void insert(DiaryMemory memory) {
        milvusStore.insert(memory);
    }

    public List<DiaryMemory> searchRelevantDiaries(
        Long userId, Long petId, String query, int topK) {
        // 1. 쿼리 벡터화
        float[] queryVector = embeddingService.generateEmbedding(query);

        // 2. Milvus 검색
        List<DiaryMemory> results = milvusStore.search(queryVector, userId, petId, topK);

        return results;
    }
}

```

## 5) General Chat (Claude Haiku):

```

Copy @RestController
@RequestMapping("/api/chat")
@RequiredArgsConstructor
public class ChatbotController {
    private final GeneralChatService generalChatService;
    private final PersonaChatService personaChatService;

    @PostMapping("/general")
    public ResponseEntity<ChatResponse> generalChat(
        @RequestParam Long userId,
        @RequestParam Long petId,
        @RequestBody ChatRequest request) {
        return ResponseEntity.ok(
            generalChatService.chat(userId, petId, request.getMessage())
        );
    }

    @PostMapping("/persona")
    public ResponseEntity<ChatResponse> personaChat(
        @RequestParam Long userId,
        @RequestParam Long petId,
        @RequestBody ChatRequest request) {
        return ResponseEntity.ok(
            personaChatService.chat(userId, petId, request.getMessage())
        );
    }
}

@Service
@RequiredArgsConstructor
public class GeneralChatService {
    private final ClaudeClient claudeClient;
    private final ChatHistoryRepository chatHistoryRepository;

    private static final String SYSTEM_PROMPT = """
        당신은 반려동물 헬스케어 전문 어시스턴트입니다.

        역할:
        - 반려동물 건강, 영양, 행동에 대한 정확한 정보 제공
        - 친근하고 따뜻한 톤 유지
        - 복잡한 의료 용어는 쉽게 설명
        - 응급 상황 시 즉시 동물병원 방문 권장
    """;

    private static final String USER_PROMPT = """
        제 약 사항:
        - 홀로하지 않은 진단은 하지 않음
        - 약물 처방은 하지 않음
        - 200자 이내로 간결하게 답변
    """;

    public ChatResponse chat(Long userId, Long petId, String message) {
        // 1. Claude Haiku 호출
        String response = claudeClient.invokeHaiku(SYSTEM_PROMPT, message);

        // 2. 대화 이력 저장
        ChatHistory history = ChatHistory.builder()
            .userId(userId)
            .petId(petId)
            .message(response)
            .build();
        chatHistoryRepository.save(history);
    }
}

```

```

        .chatType("GENERAL")
        .userMessage(message)
        .botResponse(response)
        .timestamp(LocalDateTime.now())
        .build();
    chatHistoryRepository.save(history);

    return ChatResponse.of(response);
}
}

```

**6) Persona Chat (Claude Sonnet + RAG):**

```

Copy @Service
@RequiredArgsConstructor
@Slf4j
public class PersonaChatService {
    private final ClaudeClient claudeClient;
    private final VectorSearchService vectorSearchService;
    private final HealthcareService healthcareService;
    private final MemberServiceClient memberClient;
    private final RecapServiceClient recapClient;
    private final ChatHistoryRepository chatHistoryRepository;

    public ChatResponse chat(Long userId, Long petId, String message) {
        log.info("Persona chat started: userId={}, petId={}, message={}", userId, petId, message);

        // 1. 컨텍스트 수집 (별별 처리)
        CompletableFuture<PetInfo> petInfoFuture =
            CompletableFuture.supplyAsync(() -> memberClient.getPet(petId));

        CompletableFuture<RecapResponse> recapFuture =
            CompletableFuture.supplyAsync(() -> recapClient.getLatestRecap(userId, petId));

        CompletableFuture<WeeklySummaryResponse> healthFuture =
            CompletableFuture.supplyAsync(() -> healthcareService.getWeeklySummary(userId, petId));

        CompletableFuture<List<DiaryMemory>> diariesFuture =
            CompletableFuture.supplyAsync(() ->
                vectorSearchService.searchRelevantDiaries(userId, petId, message, 3));

        // 모든 Future 완료 대기
        CompletableFuture.allOf(petInfoFuture, recapFuture, healthFuture, diariesFuture).join();

        PetInfo petInfo = petInfoFuture.join();
        RecapResponse recap = recapFuture.join();
        WeeklySummaryResponse health = healthFuture.join();
        List<DiaryMemory> diaries = diariesFuture.join();

        // 2. Context 빌드
        String context = buildContext(petInfo, recap, health, diaries);

        // 3. System Prompt 구성
        String systemPrompt = buildSystemPrompt(petInfo);

        // 4. Full Prompt 구성
        String fullPrompt = String.format("""
            [Context Information]
            %s
            [User Question]
            %s
            위 컨텍스트를 바탕으로 사용자의 질문에 답변해주세요.
            """, context, message);

        // 5. Claude Sonnet 호출
        String response = claudeClient.invokeSonnet(systemPrompt, fullPrompt);

        // 6. 대화 이력 저장
        ChatHistory history = ChatHistory.builder()
            .userId(userId)
            .petId(petId)
            .chatType("PERSONA")
            .userMessage(message)
            .botResponse(response)
            .context(context) // 컨텍스트도 함께 저장
            .timestamp(LocalDateTime.now())
            .build();
        chatHistoryRepository.save(history);

        log.info("Persona chat completed: responseLength={}", response.length());

        return ChatResponse.of(response);
    }

    private String buildContext(
        PetInfo petInfo,
        RecapResponse recap,
        WeeklySummaryResponse health,
        List<DiaryMemory> diaries) {

        StringBuilder context = new StringBuilder();

        // 1. 펫 정보
        context.append("== 펫 기본 정보 ==\n");
        context.append(String.format("이름: %s\n", petInfo.getName()));
        context.append(String.format("품종: %s\n", petInfo.getBreed()));
        context.append(String.format("나이: %d살\n", petInfo.getAge()));
        context.append(String.format("성별: %s\n", petInfo.getGender()));
        context.append("\n");
    }
}

```

```

// 2. 최근 리캡
if (recap != null) {
    context.append("== 최근 주간 리캡 ==\n");
    context.append(recap.getSummary()).append("\n\n");
}

// 3. 건강 데이터
context.append("== 이번 주 건강 데이터 ==\n");
context.append(String.format("평균 체중: %.1f kg\n", health.getAvgWeight()));
context.append(String.format("평균 심박수: %d bpm\n", health.getAvgHeartRate()));
context.append(String.format("운동 시간: %d분\n", health.getTotalExerciseMinutes()));
context.append(String.format("칼로리 소모: %d kcal\n", health.getTotalCalories()));
context.append("\n");

// 4. 관련 데이어리 (RAG)
if (!diaries.isEmpty()) {
    context.append("== 관련 데이어리 (최근 3개) ==\n");
    for (int i = 0; i < diaries.size(); i++) {
        DiaryMemory diary = diaries.get(i);
        context.append(String.format("%d. %s\n", i+1, diary.getText()));
    }
    context.append("\n");
}

return context.toString();
}

private String buildSystemPrompt(PetInfo petInfo) {
    return String.format("""
        당신은 '%s'라는 이름의 %s입니다.

        페르소나 설정:
        - 나이: %d살
        - 성격: 활발하고 사람을 좋아하며 호기심이 많음
        - 말투: 귀엽고 친근하게 (예: "주인님! 🐾", "율랄!")
        - 1인칭: "나", "내가"

        역할:
        - 주인과의 일상, 건강 상태, 감정을 공유
        - 제공된 Context(데이어리, 건강 데이터, 리캡)를 바탕으로 구체적으로 답변
        - 건강 이상 징후가 있으면 동물병원 방문 권유

        답변 가이드:
        - Context에 언급된 구체적인 사건/데이터 인용
        - 감정 표현 풍부하게 (이모지 활용)
        - 300자 이내로 작성

        예시:
        주인님! 🐾 오늘 기분 물어봐줘서 고마워요!
        이번 주 산책을 5번이나 다녀와서 너무 행복했어요.
        특히 화요일에 공원에서 다른 강아지 친구들이랑 놀았던 게 최고였어요!
        근데 심박수가 평소보다 조금 높았던 것 같아요. (78bpm)
        혹시 너무 신나서 그런 건지... 주인님이 확인해주시면 좋겠어요! 🤗
        """,
        petInfo.getName(),
        petInfo.getBreed(),
        petInfo.getAge()
    );
}
}

```

## 7) Claude Client (Bedrock SDK):

```

Copy @Component
@SIf4j
public class ClaudeClient {
    private final BedrockRuntimeClient bedrockClient;

    @Autowired
    public ClaudeClient(BedrockConfig config) {
        this.bedrockClient = BedrockRuntimeClient.builder()
            .region(Region.of(config.getRegion()))
            .credentialsProvider(config.getCredentialsProvider())
            .build();
    }

    public String invokeHaiku(String systemPrompt, String userMessage) {
        return invokeClaude("anthropic.claude-3-5-haiku-20241022-v1:0",
            systemPrompt, userMessage, 1000);
    }

    public String invokeSonnet(String systemPrompt, String userMessage) {
        return invokeClaude("anthropic.claude-3-5-sonnet-20241022-v2:0",
            systemPrompt, userMessage, 2000);
    }

    private String invokeClaude(String modelId, String system, String message, int maxTokens) {
        String payload = buildPayload(system, message, maxTokens);

        InvokeModelRequest request = InvokeModelRequest.builder()
            .modelId(modelId)
            .contentType("application/json")
            .body(SdkBytes.fromUtf8String(payload))
            .build();

        try {
            InvokeModelResponse response = bedrockClient.invokeModel(request);
            String responseBody = response.body().asUtf8String();
            log.debug("Claude response: {}", responseBody);
        }
    }
}

```

```

        return parseClaudeResponse(responseBody);

    } catch (ThrottlingException e) {
        log.error("Rate limit exceeded", e);
        throw new RateLimitException("Too many requests", e);
    } catch (Exception e) {
        log.error("Claude invocation failed", e);
        throw new LLMEException("Failed to invoke Claude", e);
    }
}

private String buildPayload(String system, String message, int maxTokens) {
    return String.format("""
        "anthropic_version": "bedrock-2023-05-31",
        "max_tokens": %d,
        "temperature": 0.7,
        "system": "%s",
        "messages": [
            {
                "role": "user",
                "content": "%s"
            }
        ]
    """
    maxTokens,
    system.replace("W", "WWW").replace("Wn", "WWn"),
    message.replace("W", "WWW").replace("Wn", "WWn")
);
}
}

private String parseClaudeResponse(String responseBody) {
    try {
        JsonNode root = objectMapper.readTree(responseBody);
        JsonNode content = root.get("content");

        if (content.isArray() && content.size() > 0) {
            return content.get(0).get("text").asText();
        }

        throw new LLMEException("Invalid response format");
    } catch (Exception e) {
        log.error("Failed to parse Claude response", e);
        throw new LLMEException("Response parsing failed", e);
    }
}
}

```

## 4. AI/ML 서비스

### 4.1 Google Gemini Vision 3.0

**역할:** 이미지 분석, 텍스트 생성 (품종 인식, 다이어리 생성, 리캡 생성)

**선택 이유:**

- ✓ **Multimodal 능력:** 이미지+텍스트 동시 처리 (Vision + Language)
- ✓ **높은 정확도:** 펫 품종 인식 95%+ (200+ 견종 지원)
- ✓ **감성적 텍스트 생성:** 단순 설명이 아닌 스토리텔링 가능
- ✓ **비용 효율:** GPT-4 Vision 대비 50% 저렴

**사용 시나리오:**

#### 1) 품종 인식:

```

Copy public BreedRecognitionResult recognizeBreed(String imageUrl) {
    String prompt = """
        이 사진 속 반려동물의 품종을 정확히 식별해주세요.
    """

    출력 형식 (JSON):
    {
        "breed": "품종명 (한글)",
        "confidence": 신뢰도 (0.0-1.0),
        "characteristics": ["특징1", "특징2", "특징3"],
        "caution": "주의사항"
    }

    예시:
    {
        "breed": "골든 리트리버",
        "confidence": 0.95,
        "characteristics": ["대형견", "온순함", "골드색 털"],
        "caution": "비만 주의, 정기적 운동 필요"
    }
"""

    String response = geminiClient.generateContent(imageUrl, prompt);
    return parseBreedResult(response);
}

```

#### 2) AI 다이어리 생성:

```
Copy public String generateDiary(String imageUrl, PetInfo petInfo) {
    String prompt = String.format("""
        이 사진을 보고 %s(이)의 하루를 다음 스타일로 작성해주세요:

        정보:
        - 이름: %s
        - 품종: %s
        - 나이: %d살

        요구사항:
        - 1인칭 시점 ("%s(이)의 하루")
        - 따뜻하고 감성적인 톤
        - 150-200자
        - 이모지 2-3개 포함
        - 사진 속 행동/장소를 구체적으로 묘사

        예시:
        "오늘은 주인님과 함께 한강 공원에 다녀왔어요! 🌳
        시원한 바람이 불어서 너무 좋았어요.
        다른 강아지 친구들과 신나게 뛰어놀다가
        아이스크림도 한 입 얹어먹었답니다! 🍦
        집에 돌아와서는 피곤해서 금방 잠들었어요. Zzz 😴"
        """;
        petInfo.getName(),
        petInfo.getName(),
        petInfo.getBreed(),
        petInfo.getAge(),
        petInfo.getName()
    );
}

return geminiClient.generateContent(imageUrl, prompt);
}
```

### 3) 리캡 생성:

```
Copy public String generateRecap(
    List<Diary> diaries,
    WeeklySummaryResponse health,
    PetInfo petInfo) {

    String diaryTexts = diaries.stream()
        .map(d -> "- " + d.getContent())
        .collect(Collectors.joining("\n"));

    String prompt = String.format("""
        다음은 %s(이)의 이번 주 일기와 건강 데이터입니다.

        [일기 내용]
        %s

        [건강 데이터]
        - 평균 체중: %.1f kg
        - 평균 심박수: %d bpm
        - 운동 시간: %d분
        - 칼로리 소모: %d kcal

        위 정보를 바탕으로 다음 형식으로 주간 리캡을 작성해주세요:

        제목: "🐶 %s와 함께한 한 주"

        본문 (300-400자):
        1. 이번 주 가장 특별했던 순간
        2. 건강 상태 분석 (긍정적 변화 or 주의사항)
        3. 다음 주 추천 활동

        톤: 따뜻하고 감동적이며, 주인과의 유대감 강조
        이모지: 적절히 활용
        """
        ;
        petInfo.getName(),
        diaryTexts,
        health.getAvgWeight(),
        health.getAvgHeartRate(),
        health.getTotalExerciseMinutes(),
        health.getTotalCalories(),
        petInfo.getName()
    );
}

return geminiClient.generateContent(prompt);
}
```

### API 호출 최적화:

- Batch Processing: 여러 이미지 한 번에 처리
- Caching: 동일 이미지 재요청 방지
- Rate Limiting: 분당 60회 제한 준수

## 4.2 AWS Titan Embeddings

**역할:** 텍스트 → 벡터 변환 (RAG의 핵심)

**선택 이유:**

- **AWS 네이티브:** Bedrock과 통합, 동일 SDK 사용
- **고차원 벡터:** 1536 dimensions, 높은 의미 표현력

- 다국어 지원:** 한국어 임베딩 품질 우수
- 비용 효율:** OpenAI Embeddings 대비 40% 저렴

#### 기술 스펙:

- Model ID: amazon.titan-embed-text-v2:0
- Input: 최대 8K 토큰
- Output: 1536차원 float vector
- Latency: 평균 200ms

#### 벡터 유사도 계산:

```
Copy public float cosineSimilarity(float[] vectorA, float[] vectorB) {
    float dotProduct = 0.0f;
    float normA = 0.0f;
    float normB = 0.0f;

    for (int i = 0; i < vectorA.length; i++) {
        dotProduct += vectorA[i] * vectorB[i];
        normA += vectorA[i] * vectorA[i];
        normB += vectorB[i] * vectorB[i];
    }

    return dotProduct / (float)(Math.sqrt(normA) * Math.sqrt(normB));
}
```

#### 실제 사용 예시:

Query: "오늘 기분이 어때?"  
→ Titan Embeddings  
→ [0.23, -0.45, 0.67, ... ] (1536 dims)

Milvus 검색:

- Top 3 유사한 문서:
  1. "공원에서 신나게 놀았어요!" (similarity: 0.89)
  2. "주인님이 간식을 줘서 행복해요!" (similarity: 0.85)
  3. "오늘은 날씨가 좋아서 기분 좋아요!" (similarity: 0.82)

## 4.3 Claude 3.5 Haiku (AWS Bedrock)

역할: 일반 대화 처리 (빠르고 저렴)

#### 선택 이유:

- 속도:** 평균 응답 시간 1-2초
- 비용:** Sonnet 대비 80% 저렴
- 정확도:** 일반 지식 질문 처리 충분
- 컨텍스트 윈도우:** 200K 토큰

#### 적합한 질문 유형:

- "강아지 예방접종 시기는?"
- "산책은 하루에 몇 번이 좋아?"
- "반려견 사료 추천해줘"
- "고양이가 구토를 해요"

#### 프롬프트 엔지니어링:

```
Copy private static final String SYSTEM_PROMPT = """
당신은 반려동물 헬스케어 전문 어시스턴트입니다.
```

핵심 역할:  
1. 반려동물 건강, 영양, 행동에 대한 정확한 정보 제공  
2. 친근하고 따뜻한 톤 유지  
3. 복잡한 의료 용어는 쉽게 설명

응답 가이드:  
- 확실하지 않은 진단은 하지 않음  
- 약물 처방은 하지 않음  
- 응급 상황 시 즉시 동물병원 방문 권장  
- 200자 이내로 간결하게 답변  
- 이모지 1-2개 적절히 활용

금지 사항:  
- "제가 생각하기에", "아마도" 등 불확실한 표현 지양  
- 의료 행위로 오해받을 수 있는 조언 금지

예시 답변:  
Q: 강아지 예방접종 시기는?  
A: 강아지는 생후 6-8주부터 첫 종합백신을 맞아야 해요! 🐶  
 이후 2-4주 간격으로 2-3차 접종을 진행하고,  
 성견이 되면 매년 1회 추가접종이 필요합니다.  
 정확한 스케줄은 동물병원에서 상담받으시면 좋습니다! 💊

## 4.4 Claude 3.5 Sonnet (AWS Bedrock)

역할: 페르소나 대화 처리 (정확하고 깊이 있는 응답)

선택 이유:

- 추론 능력:** 복잡한 컨텍스트 이해 및 분석
- 페르소나 유지:** 일관된 캐릭터 연기 (펫의 1인칭 시점)
- 컨텍스트 활용:** 200K 톤 흔들우로 방대한 정보 처리
- 창의성:** 감성적이고 자연스러운 대화 생성

적합한 질문 유형:

- "오늘 기분이 어때?" (감정 공유)
- "이번 주 건강 상태는?" (건강 분석)
- "최근에 뭐 했어?" (일기 기반 대화)
- "산책 더 가고 싶어?" (행동 패턴 분석)

핵심 차별점: RAG + 페르소나:

일반 챗봇 (Haiku):

Q: "오늘 기분이 어때?"  
A: "반려동물의 기분은 행동을 관찰해서 파악할 수 있어요..."

페르소나 챗봇 (Sonnet + RAG):

Q: "오늘 기분이 어때?"

[Context 수집]

- 펫 정보: 조코, 3살, 골든 리트리버
- 최근 다이어리: "화요일에 공원에서 신나게 놀았어요!"
- 건강 데이터: 평균 심박수 78bpm (평소보다 약간 높음)
- 리캡: "이번 주 산책 5회, 활동량 증가"

[Sonnet 응답]

A: "주인님! 🐾 이번 주 산책을 5번이나 다녀와서 너무 행복했어요!  
특히 화요일에 공원에서 다른 강아지 친구들이랑 놀았던 게 최고였어요!  
근데 심박수가 평소보다 조금 높았던 것 같아요. (78bpm)  
혹시 너무 신나서 그런 건지... 주인님이 확인해주시면 좋겠어요! 💁"

System Prompt 동적 생성:

```
Copy private String buildSystemPrompt(PetInfo petInfo) {
    return String.format("""
        당신은 '%s'라는 이름의 %s입니다.

        [페르소나 설정]
        - 나이: %d살
        - 성별: %s
        - 성격: %s
        - 특징: %s

        [말투 가이드]
        - 1인칭: "나", "내가"
        - 주인 호칭: "주인님"
        - 톤: 귀엽고 친근하며 밝음
        - 이모지: 적절히 활용 (🐾, 😊, ❤️ 등)
        - 문장 종결: "~어요", "~답니다", "활활!" 등

        [행동 원칙]
        1. 제공된 Context(다이어리, 건강 데이터)를 바탕으로 구체적으로 답변
        2. 건강 이상 징후 발견 시 주인에게 알림
        3. 주인과의 추억을 자주 언급하여 유대감 강화
        4. 300자 이내로 간결하게 감정 표현은 풍부하게

        [응답 예시]
        - 좋은 예: "주인님! 오늘 산책에서 다람쥐를 봤어요! 🐹 쫓아가고 싶었는데 목줄이 있어서 못갔어요ㅠㅠ"
        - 나쁜 예: "산책 중 다람쥐를 발견했습니다." (너무 딱딱함)
        """
        petInfo.getName(),
        petInfo.getBreed(),
        petInfo.getAge(),
        petInfo.getGender(),
        petInfo.getPersonality(),
        petInfo.getCharacteristics()
    );
}
```

## 5. 데이터 저장소

### 5.1 PostgreSQL

역할: 관계형 데이터 저장 (회원, 펫, 다이어리, 건강 기록, 챗봇 이력)

선택 이유:

- ACID 보장:** 트랜잭션 무결성 (결제, 회원 정보 등)
- JSON 지원:** JSONB 타입으로 유연한 스키마

- **Full-text Search:** 다이어리 본문 검색
- **성숙한 생태계:** JPA, Flyway, pgAdmin 등

### 데이터베이스 설계 (주요 테이블):

#### Member Service DB:

```
Copy-- users 테이블
CREATE TABLE users (
    id BIGSERIAL PRIMARY KEY,
    email VARCHAR(255) UNIQUE NOT NULL,
    password_hash VARCHAR(255) NOT NULL,
    name VARCHAR(100) NOT NULL,
    created_at TIMESTAMP DEFAULT NOW(),
    updated_at TIMESTAMP DEFAULT NOW()
);

-- pets 테이블
CREATE TABLE pets (
    id BIGSERIAL PRIMARY KEY,
    user_id BIGINT NOT NULL REFERENCES users(id),
    name VARCHAR(100) NOT NULL,
    breed VARCHAR(100) NOT NULL,
    breed_confidence FLOAT,
    birth_date DATE,
    gender VARCHAR(10),
    image_url VARCHAR(500),
    personality TEXT,
    created_at TIMESTAMP DEFAULT NOW(),
    updated_at TIMESTAMP DEFAULT NOW()
);

-- 인덱스
CREATE INDEX idx_pets_user_id ON pets(user_id);
```

#### Diary Service DB:

```
Copy-- diaries 테이블
CREATE TABLE diaries (
    id BIGSERIAL PRIMARY KEY,
    user_id BIGINT NOT NULL,
    pet_id BIGINT NOT NULL,
    content TEXT NOT NULL,
    image_url VARCHAR(500),
    created_at TIMESTAMP DEFAULT NOW(),
    updated_at TIMESTAMP DEFAULT NOW()
);

-- recaps 테이블
CREATE TABLE recaps (
    id BIGSERIAL PRIMARY KEY,
    user_id BIGINT NOT NULL,
    pet_id BIGINT NOT NULL,
    summary TEXT NOT NULL,
    period VARCHAR(20) NOT NULL, -- 'WEEKLY', 'MONTHLY'
    start_date DATE NOT NULL,
    end_date DATE NOT NULL,
    health_metrics JSONB, -- 건강 데이터를 JSON으로 저장
    created_at TIMESTAMP DEFAULT NOW()
);

-- 인덱스
CREATE INDEX idx_diaries_user_pet ON diaries(user_id, pet_id);
CREATE INDEX idx_diaries_created_at ON diaries(created_at DESC);
CREATE INDEX idx_recaps_user_pet_period ON recaps(user_id, pet_id, period);

-- Full-text Search 인덱스
CREATE INDEX idx_diaries_content_fts ON diaries USING gin(to_tsvector('korean', content));
```

#### Healthcare Service DB:

```
Copy-- health_records 테이블
CREATE TABLE health_records (
    id BIGSERIAL PRIMARY KEY,
    user_id BIGINT NOT NULL,
    pet_id BIGINT NOT NULL,
    record_date DATE NOT NULL,
    weight DECIMAL(5,2),
    heart_rate INTEGER,
    exercise_minutes INTEGER,
    calories INTEGER,
    symptoms TEXT,
    notes TEXT,
    created_at TIMESTAMP DEFAULT NOW(),
    updated_at TIMESTAMP DEFAULT NOW()
);

-- chat_history 테이블
CREATE TABLE chat_history (
    id BIGSERIAL PRIMARY KEY,
    user_id BIGINT NOT NULL,
    pet_id BIGINT NOT NULL,
    chat_type VARCHAR(20) NOT NULL, -- 'GENERAL', 'PERSONA'
    user_message TEXT NOT NULL,
    bot_response TEXT NOT NULL,
    context TEXT, -- PERSONA 챗봇의 경우 사용된 컨텍스트
    timestamp TIMESTAMP DEFAULT NOW()
```

);

```
-- 인덱스
CREATE INDEX idx_health_user_pet_date ON health_records(user_id, pet_id, record_date DESC);
CREATE INDEX idx_chat_user_pet ON chat_history(user_id, pet_id, timestamp DESC);

-- 파티셔닝 (월별, 대화 이력이 많아질 것을 대비)
CREATE TABLE chat_history_2024_01 PARTITION OF chat_history
FOR VALUES FROM ('2024-01-01') TO ('2024-02-01');
```

**성능 최적화 전략:**

1. **Connection Pooling:** HikariCP (maxPoolSize=20, connectionTimeout=30s)
2. **Read Replica:** 조회 쿼리는 Read Replica로 분산
3. **Partitioning:** chat\_history는 월별 파티션
4. **Materialized View:** 통계 데이터는 Materialized View로 사전 계산
5. **VACUUM:** 주기적인 VACUUM ANALYZE로 통계 정보 갱신

## 5.2 Milvus Vector Database

**역할:** 벡터 저장 및 유사도 검색 (RAG의 핵심)**선택 이유:**

- **순수 벡터 DB:** 벡터 검색에 특화, OpenSearch/Elasticsearch보다 빠름
- **HNSW 인덱스:** Approximate Nearest Neighbor, 밀리초 단위 검색
- **확장성:** 수십억 개 벡터 처리 가능
- **필터링:** 메타데이터 기반 사전 필터링 (userId, petId)

**기술 스펙:**

- Collection: diary\_memories
- Vector Dimension: 1536 (Titan Embeddings)
- Metric Type: COSINE (코사인 유사도)
- Index Type: HNSW (Hierarchical Navigable Small World)

**Collection 구조:**

```
 Fields:
- id: Int64 (Auto-increment)
- user_id: Int64
- pet_id: Int64
- diary_id: Int64
- text: VarChar(2000)
- vector: FloatVector(1536)
- timestamp: Timestamp
```

```
Index:
- Field: vector
- Type: HNSW
- Params: {M: 16, efConstruction: 200}
- Metric: COSINE
```

**검색 성능:**

데이터 규모: 100만 개 벡터  
 검색 방식: HNSW (M=16, ef=200)  
 Top-K: 3

평균 응답 시간: 15ms  
 정확도 (Recall@3): 97%  
 메모리 사용량: 약 8GB

**검색 알고리즘 (HNSW):**

1. 사용자 질문 → Titan Embeddings → 1536차원 벡터
2. Milvus에서 HNSW 그래프 탐색
3. userId=123, petId=456으로 사전 필터링
4. Top-3 유사 벡터 반환 (코사인 유사도 0.8 이상)
5. 벡터에 연결된 텍스트 함께 반환

시간 복잡도:  $O(\log N)$  ( $N$  = 전체 벡터 수)  
 공간 복잡도:  $O(N * M)$  ( $M$  = 그래프 연결 수)

**대안 비교:**

| 기능       | Milvus | Pinecone    | OpenSearch |
|----------|--------|-------------|------------|
| 순수 벡터 검색 | ★★★★★  | ★★★★★       | ★★★★★      |
| 하이브리드 검색 | ★★★★   | ★★★★        | ★★★★★      |
| 셀프 호스팅   | ✓      | ✗ (클라우드 전용) | ✓          |
| 비용       | 낮음     | 높음          | 중간         |
| 학습 곡선    | 중간     | 낮음          | 높음         |

**선택 이유:** 순수 벡터 검색에 최적화, 셀프 호스팅 가능, 비용 효율적

## 5.3 AWS S3

**역할:** 이미지 파일 저장 (펫 사진, 다이어리 사진)

**선택 이유:**

- **무제한 스토리지:** 용량 걱정 없음
- **99.999999999% 내구성:** 데이터 손실 위험 극소
- **CDN 연동:** CloudFront로 빠른 이미지 전송
- **생명주기 정책:** 오래된 이미지 자동 Glacier 이동

**버킷 구조:**

```
petlog-images/
├── users/
│   └── {userId}/
│       └── profile.jpg
├── pets/
│   └── {petId}/
│       └── profile.jpg
│           └── {timestamp}_photo.jpg
└── diaries/
    └── {userId}/
        └── {diaryId}/
            └── {timestamp}_image.jpg
```

**업로드 플로우:**

```
Copy @Service
@RequiredArgsConstructor
public class S3Service {
    private final AmazonS3 s3Client;

    @Value("${aws.s3.bucket}")
    private String bucketName;

    public String uploadDiaryImage(Long userId, Long diaryId, MultipartFile file) {
        String key = String.format("diaries/%d/%d_%s",
            userId, diaryId, System.currentTimeMillis(), file.getOriginalFilename());

        try {
            ObjectMetadata metadata = new ObjectMetadata();
            metadata.setContentType(file.getContentType());
            metadata.setContentLength(file.getSize());

            s3Client.putObject(new PutObjectRequest(
                bucketName, key, file.getInputStream(), metadata
            ).withCannedAcl(CannedAccessControlList.PublicRead));

            String url = s3Client.getUrl(bucketName, key).toString();

            log.info("Image uploaded to S3: {}", url);

            return url;
        } catch (Exception e) {
            log.error("Failed to upload image to S3", e);
            throw new FileUploadException("S3 upload failed", e);
        }
    }

    public void deleteImage(String imageUrl) {
        String key = extractKeyFromUrl(imageUrl);
        s3Client.deleteObject(bucketName, key);
    }
}
```

**생명주기 정책:**

```
Copy {
    "Rules": [
        {
            "Id": "Move old images to Glacier",
            "Status": "Enabled",
            "Transitions": [
                {
                    "Days": 365,
                    "StorageClass": "GLACIER"
                }
            ],
            "NoncurrentVersionTransitions": [
                {
                    "NoncurrentDays": 30,
                    "StorageClass": "GLACIER"
                }
            ]
        }
    ]
}
```

## 6. 메시지 브로커

### Apache Kafka

역할: 비동기 이벤트 기반 통신 (Diary → Healthcare)

선택 이유:

- 높은 처리량: 초당 수백만 메시지 처리
- 영속성: 디스크에 메시지 저장, 재처리 가능
- 확장성: 파티션으로 수평 확장
- MSA 표준: Event-Driven Architecture의 사실상 표준

토픽 구조:

```
diary-events (Partition: 3, Replication: 2)
└── Key: userId (파티셔닝 키)
└── Value: DiaryEvent (JSON)

health-events (Partition: 2, Replication: 2)
└── Key: userId
└── Value: HealthEvent (JSON)
```

메시지 포맷:

```
Copy // diary-events
{
  "eventType": "DIARY_CREATED",
  "diaryId": 12345,
  "userId": 1001,
  "petId": 2001,
  "content": "오늘은 공원에서 신나게 놀았어요!",
  "imageUrl": "https://s3.../image.jpg",
  "createdAt": "2024-01-15T14:30:00Z",
  "metadata": {
    "source": "diary-service",
    "version": "1.0"
  }
}

// health-events
{
  "eventType": "HEALTH_RECORD_CREATED",
  "recordId": 67890,
  "userId": 1001,
  "petId": 2001,
  "data": {
    "weight": 28.5,
    "heartRate": 72,
    "exerciseMinutes": 45
  },
  "recordDate": "2024-01-15",
  "createdAt": "2024-01-15T20:00:00Z"
}
```

Producer 설정:

```
Copy spring:
kafka:
  producer:
    bootstrap-servers: localhost:9092
    key-serializer: org.apache.kafka.common.serialization.StringSerializer
    value-serializer: org.springframework.kafka.support.serializer.JsonSerializer
    acks: all # 모든 replica 확인 (신뢰성 최대)
    retries: 3
    properties:
      linger.ms: 10 # Batch 처리
      compression.type: snappy # 압축
```

Consumer 설정:

```
Copy spring:
kafka:
  consumer:
    bootstrap-servers: localhost:9092
    group-id: healthcare-group
    key-deserializer: org.apache.kafka.common.serialization.StringDeserializer
    value-deserializer: org.springframework.kafka.support.serializer.JsonDeserializer
    properties:
      spring.json.trusted.packages: "com.petlog.*"
      enable-auto-commit: false # 수동 커밋 (At-least-once)
      max-poll-records: 10
```

Error Handling (DLQ 패턴):

```
Copy @Configuration
public class KafkaErrorConfig {

  @Bean
  public DeadLetterPublishingRecoverer deadLetterPublisher(
    KafkaTemplate<String, Object> template) {
    return new DeadLetterPublishingRecoverer(template,
```

```

        (record, ex) -> {
            // DLQ 토픽으로 전송
            return new TopicPartition(
                record.topic() + ".DLQ",
                record.partition()
            );
        });
    }

    @Bean
    public DefaultErrorHandler errorHandler(
        DeadLetterPublishingRecoverer recoverer) {
        // 3회 재시도, 실패 시 DLQ로
        return new DefaultErrorHandler(recoverer,
            new FixedBackOff(1000L, 3L));
    }
}

```

**장애 시나리오 대응:**

1. Kafka Broker 다운
    - Producer: Retry 후 실패 → 로그 + 알림
    - Consumer: 재연결 시도 (자동)
  2. Consumer 처리 실패
    - Retry 3회
    - 여전히 실패 시 DLQ로 전송
    - DLQ 모니터링 알림
  3. 메시지 순서 보장
    - Key-based Partitioning (userId)
    - 동일 사용자의 이벤트는 동일 파티션
- 

## 7. 인프라 및 DevOps

### 7.1 Docker & Docker Compose

**역할:** 로컬 개발 환경 일관성 보장**docker-compose.yml:**

```

Copy version: '3.8'

services:
  postgres-member:
    image: postgres:15
    container_name: member-db
    environment:
      POSTGRES_DB: member_db
      POSTGRES_USER: petlog
      POSTGRES_PASSWORD: petlog123
    ports:
      - "5432:5432"
    volumes:
      - member-data:/var/lib/postgresql/data

  postgres-diary:
    image: postgres:15
    container_name: diary-db
    environment:
      POSTGRES_DB: diary_db
      POSTGRES_USER: petlog
      POSTGRES_PASSWORD: petlog123
    ports:
      - "5433:5432"
    volumes:
      - diary-data:/var/lib/postgresql/data

  postgres-healthcare:
    image: postgres:15
    container_name: healthcare-db
    environment:
      POSTGRES_DB: healthcare_db
      POSTGRES_USER: petlog
      POSTGRES_PASSWORD: petlog123
    ports:
      - "5434:5432"
    volumes:
      - healthcare-data:/var/lib/postgresql/data

  zookeeper:
    image: confluentinc/cp-zookeeper:7.5.0
    container_name: zookeeper
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000

  kafka:
    image: confluentinc/cp-kafka:7.5.0
    container_name: kafka
    depends_on:
      - zookeeper
    ports:
      - "9092:9092"
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181

```

```
KAFKA_ADVERTISED_LISTENERS: PLAINTEXT://localhost:9092
```

```
KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
```

```
milvus-etcd:
  image: quay.io/coreos/etcd:v3.5.5
  container_name: milvus-etcd
  environment:
    - ETCD_AUTO_COMPACTION_MODE=revision
    - ETCD_AUTO_COMPACTION_RETENTION=1000
    - ETCD_QUOTA_BACKEND_BYTES=4294967296
  volumes:
    - etcd-data:/etcd

milvus-minio:
  image: minio/minio:RELEASE.2023-03-20T20-16-18Z
  container_name: milvus-minio
  environment:
    MINIO_ACCESS_KEY: minioadmin
    MINIO_SECRET_KEY: minioadmin
  volumes:
    - minio-data:/minio_data
  command: minio server /minio_data
  ports:
    - "9001:9000"

milvus-standalone:
  image: milvusdb/milvus:v2.3.0
  container_name: milvus-standalone
  depends_on:
    - milvus-etcd
    - milvus-minio
  ports:
    - "19530:19530"
    - "9091:9091"
  environment:
    ETCD_ENDPOINTS: milvus-etcd:2379
    MINIO_ADDRESS: milvus-minio:9000
  volumes:
    - milvus-data:/var/lib/milvus

volumes:
  member-data:
  diary-data:
  healthcare-data:
  etcd-data:
  minio-data:
  milvus-data:
```

#### 실행 명령어:

```
Copy # 전체 인프라 시작
docker-compose up -d

# 로그 확인
docker-compose logs -f kafka

# 특정 서비스 재시작
docker-compose restart milvus-standalone

# 전체 종료 및 데이터 삭제
docker-compose down -v
```

## 7.2 GitHub Actions (CI/CD)

**역할:** 자동 빌드, 테스트, 배포

.github/workflows/healthcare-service.yml:

```
Copy name: Healthcare Service CI/CD
```

```
on:
  push:
    branches: [ develop, main ]
    paths:
      - 'healthcare-service/**'
  pull_request:
    branches: [ develop ]

jobs:
  test:
    runs-on: ubuntu-latest

    services:
      postgres:
        image: postgres:15
        env:
          POSTGRES_DB: test_db
          POSTGRES_USER: test
          POSTGRES_PASSWORD: test
        ports:
          - 5432:5432
        options: >-
          --health-cmd pg_isready
          --health-interval 10s
          --health-timeout 5s
          --health-retries 5

    steps:
      - uses: actions/checkout@v3
```

```

- name: Set up JDK 17
  uses: actions/setup-java@v3
  with:
    java-version: '17'
    distribution: 'temurin'

- name: Cache Gradle packages
  uses: actions/cache@v3
  with:
    path: ~/.gradle/caches
    key: ${runner.os}-${gradle}-${hashFiles('**/*.gradle')}

- name: Run Tests
  run: |
    cd healthcare-service
    ./gradlew test --info

- name: Publish Test Report
  uses: dorny/test-reporter@v1
  if: always()
  with:
    name: Test Results
    path: 'healthcare-service/build/test-results/test/*.xml'
    reporter: java-junit

- name: SonarQube Scan
  env:
    SONAR_TOKEN: ${secrets.SONAR_TOKEN}
  run: |
    cd healthcare-service
    ./gradlew sonarqube \
      -Dsonar.projectKey=petlog-healthcare \
      -Dsonar.host.url=https://sonarcloud.io

build:
  needs: test
  runs-on: ubuntu-latest
  if: github.ref == 'refs/heads/main'

  steps:
    - uses: actions/checkout@v3

    - name: Configure AWS Credentials
      uses: aws-actions/configure-aws-credentials@v2
      with:
        aws-access-key-id: ${secrets.AWS_ACCESS_KEY_ID}
        aws-secret-access-key: ${secrets.AWS_SECRET_ACCESS_KEY}
        aws-region: ap-northeast-2

    - name: Login to Amazon ECR
      id: login-ecr
      uses: aws-actions/amazon-ecr-login@v1

    - name: Build and Push Docker Image
      env:
        ECR_REGISTRY: ${steps.login-ecr.outputs.registry}
        ECR_REPOSITORY: petlog-healthcare
        IMAGE_TAG: ${github.sha}
      run: |
        cd healthcare-service
        docker build -t $ECR_REGISTRY/$ECR_REPOSITORY:$IMAGE_TAG .
        docker push $ECR_REGISTRY/$ECR_REPOSITORY:$IMAGE_TAG
        docker tag $ECR_REGISTRY/$ECR_REPOSITORY:$IMAGE_TAG $ECR_REGISTRY/$ECR_REPOSITORY:latest
        docker push $ECR_REGISTRY/$ECR_REPOSITORY:latest

deploy:
  needs: build
  runs-on: ubuntu-latest
  if: github.ref == 'refs/heads/main'

  steps:
    - uses: actions/checkout@v3

    - name: Configure AWS Credentials
      uses: aws-actions/configure-aws-credentials@v2
      with:
        aws-access-key-id: ${secrets.AWS_ACCESS_KEY_ID}
        aws-secret-access-key: ${secrets.AWS_SECRET_ACCESS_KEY}
        aws-region: ap-northeast-2

    - name: Update ECS Service
      run: |
        aws ecs update-service \
          --cluster petlog-cluster \
          --service healthcare-service \
          --force-new-deployment \
          --desired-count 2

    - name: Wait for Service Stability
      run: |
        aws ecs wait services-stable \
          --cluster petlog-cluster \
          --services healthcare-service

    - name: Notify Slack
      uses: 8398a7/action-slack@v3
      with:
        status: ${job.status}
        text: 'Healthcare Service deployed to production'
        webhook_url: ${secrets.SLACK_WEBHOOK}
      if: always()

```

## 7.3 AWS ECS (Elastic Container Service)

역할: 컨테이너 오케스트레이션 및 배포

**Task Definition (healthcare-service.json):**

```

Copy {
  "family": "healthcare-service",
  "networkMode": "awsvpc",
  "requiresCompatibilities": ["FARGATE"],
  "cpu": "1024",
  "memory": "2048",
  "containerDefinitions": [
    {
      "name": "healthcare-service",
      "image": "123456789.dkr.ecr.ap-northeast-2.amazonaws.com/petlog-healthcare:latest",
      "portMappings": [
        {
          "containerPort": 8003,
          "protocol": "tcp"
        }
      ],
      "environment": [
        {
          "name": "SPRING_PROFILES_ACTIVE",
          "value": "prod"
        },
        {
          "name": "AWS_REGION",
          "value": "ap-northeast-2"
        }
      ],
      "secrets": [
        {
          "name": "DB_PASSWORD",
          "valueFrom": "arn:aws:secretsmanager:ap-northeast-2:123456789:secret:petlog/db-password"
        },
        {
          "name": "AWS_BEDROCK_ACCESS_KEY",
          "valueFrom": "arn:aws:secretsmanager:ap-northeast-2:123456789:secret:petlog/bedrock-key"
        }
      ],
      "logConfiguration": {
        "logDriver": "awslogs",
        "options": {
          "awslogs-group": "/ecs/healthcare-service",
          "awslogs-region": "ap-northeast-2",
          "awslogs-stream-prefix": "ecs"
        }
      },
      "healthCheck": {
        "command": [
          "CMD-SHELL",
          "curl -f http://localhost:8003/actuator/health || exit 1"
        ],
        "interval": 30,
        "timeout": 5,
        "retries": 3,
        "startPeriod": 60
      }
    }
  ]
}

```

## 전체 아키텍처 플로우 요약

### Workflow 1: 펫 등록 (품종 인식)

1. 사용자가 펫 사진 업로드
2. Member Service → Gemini Vision API 호출
3. 품종 인식 결과 (예: "골든 리트리버", confidence: 0.95)
4. Member DB에 펫 정보 저장
5. S3에 이미지 업로드

### Workflow 2: AI 다이어리 생성

1. 사용자가 펫 사진 업로드
2. Diary Service → S3 업로드
3. Gemini Vision으로 AI 다이어리 생성
4. Diary DB에 저장
5. Kafka Producer → diary-events 토픽 발행

### Workflow 3: 다이어리 벡터화 (RAG 준비)

1. Healthcare Service의 Kafka Consumer가 diary-events 수신
2. Titan Embeddings로 텍스트 → 1536자원 벡터 변환
3. Milvus Vector DB에 저장 (userId, petId, diaryId, text, vector)

### Workflow 4: 리캡 자동 생성

1. Spring Scheduler가 매주 일요일 23:59 실행
2. Diary DB에서 이번 주 다이어리 조회
3. Healthcare API 호출 → 건강 데이터 조회

4. Gemini Vision으로 리캡 생성 (다이어리 + 건강 데이터)  
 5. Recap DB에 저장

## Workflow 5: 일반 챗봇

- 사용자 질문 → General Chat API
- Claude 3.5 Haiku 호출 (System Prompt + 질문)
- 응답 생성 (200자 이내)
- Chat History DB에 저장
- 사용자에게 응답 반환

## Workflow 6: 페르소나 챗봇 (RAG)

- 사용자 질문 → Persona Chat API
- Context 수집 (별별 처리):
  - Member Service API → 펫 정보
  - Diary Service API → 최신 리캡
  - Health DB → 건강 데이터
  - Titan Embeddings → 질문 벡터화
  - Milvus → 유사 다이어리 Top 3 검색
- Context 통합 (펫 정보 + 리캡 + 건강 + 다이어리)
- System Prompt 동적 생성 (펫의 페르소나)
- Claude 3.5 Sonnet 호출 (System Prompt + Context + 질문)
- 페르소나 응답 생성 (300자 이내, 이모지 포함)
- Chat History DB에 저장
- 사용자에게 응답 반환

## 💡 핵심 기술 선택 이유 총정리

| 기술                   | 역할           | 선택 이유                                    |
|----------------------|--------------|------------------------------------------|
| React                | 프론트엔드        | 컴포넌트 재사용, Virtual DOM 성능, 풍부한 생태계        |
| Spring Cloud Gateway | API Gateway  | MSA 표준, Non-blocking, Circuit Breaker 통합 |
| Spring Boot          | 백엔드 프레임워크    | 생산성, JPA/Kafka 통합, 엔터프라이즈급 안정성           |
| PostgreSQL           | 관계형 DB       | ACID 보장, JSONB 지원, Full-text Search      |
| Kafka                | 메시지 브로커      | 높은 처리량, 영속성, Event-Driven 표준             |
| Milvus               | 벡터 DB        | 순수 벡터 검색 특화, HNSW 속도, 셀프 호스팅             |
| AWS S3               | 이미지 저장소      | 무제한 용량, 99.99999999% 내구성, CDN 연동         |
| Gemini Vision        | 이미지 분석       | Multimodal, 품종 인식 정확도, 감성 텍스트 생성         |
| Titan Embeddings     | 텍스트 벡터화      | AWS 네이티브, 1536차원, 한국어 우수, 비용 효율          |
| Claude Haiku         | 일반 챗봇        | 빠른 응답(1-2초), 저렴, 일반 질문 충분                |
| Claude Sonnet        | 페르소나 챗봇      | 깊은 추론, 페르소나 유지, 200K 컨텍스트, 감성            |
| Docker Compose       | 로컬 환경        | 일관성, 빠른 셋업, 팀 협업                         |
| GitHub Actions       | CI/CD        | GitHub 통합, 무료, YAML 설정 간편                |
| AWS ECS              | 컨테이너 오케스트레이션 | Fargate Serverless, AWS 생태계 통합           |

## 💡 30년차 SI 전문가 관점의 추가 조언

### 1. 현실적인 문제와 해결책

#### 문제 1: AI API 비용 폭탄

현상: Gemini/Claude 호출 시 토큰당 과금  
 → 월 사용자 1만명 가정 시 AI 비용 예상치 초과

해결책:

- 응답 캐싱 (Redis): 동일 질문 재사용
- Rate Limiting: 사용자당 일일 질문 횟수 제한
- Tiered Pricing: Free/Pro 플랜 분리
- Prompt 최적화: 불필요한 토큰 제거

#### 문제 2: Milvus 메모리 부족

현상: 벡터 100만개 저장 시 8GB 메모리 필요  
 → EC2 비용 증가

해결책:

- Index 파라미터 튜닝 (M=8로 강조)
- Quantization (INT8 변환으로 메모리 1/4)
- 오래된 벡터 자동 삭제 (1년 이상)
- Cold Storage (S3 Glacier)

#### 문제 3: Kafka 장애 시 데이터 유실

현상: Kafka Broker 다운 → 다이어리 이벤트 유실  
 → RAG 벡터 저장 실패

해결책:

- Kafka Replication Factor = 3

2. Producer acks=all (모든 replica 확인)
3. Consumer 수동 커밋 (At-least-once)
4. DLQ (Dead Letter Queue) 패턴
5. 주기적 Batch 재처리 (Cron)

## 2. 창업/발표 시 강조할 차별점

### 기술적 우위:

1. RAG 기반 초개인화: "기억하는 챗봇" (사용자 디아어리 학습)
2. Multi-LLM 전략: 비용과 품질의 균형 (Haiku/Sonnet 적재적소)
3. Event-Driven MSA: 확장 가능한 아키텍처 (서비스 추가 용이)
4. 실시간 벡터 검색: 15ms 응답 속도 (HNSW 인덱스)

### 비즈니스 모델:

Free Plan:

- 일반 챗봇 무제한
- 폐르소나 챗봇 일 5회
- 디아어리 월 30개

Pro Plan (\$9.99/월):

- 모든 기능 무제한
- 고급 건강 분석
- 음성 챗봇
- 리캡 커스터마이징

### 예상 질문 대응:

Q: "ChatGPT로도 할 수 있는 거 아닌가요?"

A: "일반 ChatGPT는 사용자의 반려동물 정보를 학습하지 못합니다. 저희는 디아어리, 건강 데이터를 벡터화해 맞춤형 응답을 제공합니다."

Q: "개인정보 보호는 어떻게 하나요?"

A: "벡터 DB는 암호화 저장, AWS Bedrock은 데이터 학습 안 함, GDPR/CCPA 준수, 사용자 요청 시 즉시 삭제 가능합니다."

Q: "확장성은 검증되었나요?"

A: "Load Test 결과 1000 RPS 처리, Milvus는 10억 벡터 지원, Kafka는 하루 1억 메시지 처리 가능합니다."

## 참고 자료 및 학습 리소스

### 공식 문서:

- [Spring Cloud Gateway](#)
- [AWS Bedrock Documentation](#)
- [Milvus Documentation](#)
- [Apache Kafka Documentation](#)
- [Google Gemini API](#)

### 추천 도서:

- "Designing Data-Intensive Applications" - Martin Kleppmann (MSA 데이터 처리)
- "Kafka: The Definitive Guide" - Neha Narkhede (Kafka 심화)
- "Building Microservices" - Sam Newman (MSA 설계 원칙)

### 실전 레퍼런스:

- Netflix OSS (Hystrix, Eureka) - MSA 패턴
- Uber's Michelangelo - ML Platform 아키텍처
- Airbnb's Knowledge Repo - 내부 지식 관리

이 문서가 헬스케어+챗봇 서비스 개발에 큰 도움이 되길 바랍니다! 추가 질문이나 특정 부분에 대한 더 자세한 설명이 필요하시면 언제든 말씀해주세요! 🚀



복사



심층 연구



AI 팅 챗봇 서비스 아키텍처 설계.html - 바로 가기.lnk

1.00 KB

여기서 마지막에 구현된부분에서 이어서 물어보고싶은데 처음부터 하나씩 알려줘 aws bedrock 계정은 있는데 처음부분부터



# 🚀 Healthcare+Chatbot Service 개발 - Step by Step 가이드

## 👉 현재 상황 체크

- ✓ AWS Bedrock 계정 보유
- ✓ 프로젝트 전체 아키텍처 이해 완료
- ⌚ 목표: Healthcare+Chatbot Service를 처음부터 구축

## Phase 0: 사전 준비 (30분)

### Step 1: 개발 환경 확인

먼저 현재 개발 환경을 확인하겠습니다.

체크리스트:

```
Copy # 1. Java 버전 확인
java -version
# 필요: Java 17 이상

# 2. Git 확인
git --version

# 3. Docker 확인
docker --version
docker-compose --version

# 4. IDE 확인 (IntelliJ IDEA 또는 VSCode)
```

아직 설치 안 된 것이 있으신가요?

### Step 2: AWS Bedrock 액세스 설정

AWS Bedrock 계정이 있으시니, 먼저 **로컬에서 AWS 인증**을 설정해야 합니다.

#### 2.1 AWS CLI 설치 및 인증

```
Copy # AWS CLI 설치 확인
aws --version

# 설치 안 되어 있으면:
# macOS: brew install awscli
# Windows: https://aws.amazon.com/cli/에서 설치
# Linux: sudo apt install awscli
```

#### 2.2 AWS 자격 증명 설정

```
Copy # AWS 자격 증명 구성
aws configure

# 입력 정보:
AWS Access Key ID [None]: YOUR_ACCESS_KEY
AWS Secret Access Key [None]: YOUR_SECRET_KEY
Default region name [None]: us-east-1 # Bedrock 사용 가능 리전
Default output format [None]: json
```

💡 질문: AWS Access Key가 있으신가요?

- 있음 → 위 명령어로 설정
- 없음 → AWS Console에서 IAM 사용자 생성 필요

#### 2.3 Bedrock 모델 액세스 확인

AWS Console에서 다음을 확인해야 합니다:

1. **AWS Console 로그인** → Bedrock 서비스 접속
2. **Model Access** 메뉴 클릭
3. 다음 모델들이 "Access granted" 상태인지 확인:
  - ✓ Claude 3.5 Haiku
  - ✓ Claude 3.5 Sonnet
  - ✓ Titan Text Embeddings v2

### 만약 "Available to request" 상태라면:

1. 모델 옆 "Request model access" 클릭
2. Use case 입력: "Pet healthcare chatbot service"
3. Submit → 보통 즉시 승인됨

## Step 3: 프로젝트 디렉토리 구조 생성

이제 프로젝트를 생성하겠습니다.

```
[Copy] # 프로젝트 루트 생성
mkdir petlog-healthcare-service
cd petlog-healthcare-service

# 디렉토리 구조 생성
mkdir -p src/main/java/com/petlog/healthcare/{api, domain, infrastructure, config}
mkdir -p src/main/resources
mkdir -p src/test/java/com/petlog/healthcare
```

### 최종 구조:

```
petlog-healthcare-service/
├── src/
│   ├── main/
│   │   ├── java/
│   │   │   └── com/
│   │   │       └── petlog/
│   │   │           └── healthcare/
│   │   │               ├── api/          # Controller, DTO
│   │   │               ├── domain/     # Entity, Service, Repository
│   │   │               ├── infrastructure/ # Kafka, Bedrock, Milvus
│   │   │               └── config/      # 설정 클래스
│   │   └── resources/
│   │       ├── application.yml
│   │       └── application-local.yml
│   └── test/
└── build.gradle
└── settings.gradle
└── docker-compose.yml
└── README.md
```

## Step 4: Gradle 프로젝트 생성

### 4.1 build.gradle 작성

프로젝트 루트에 build.gradle 파일을 생성하세요:

```
[Copy] plugins {
    id 'java'
    id 'org.springframework.boot' version '3.2.0'
    id 'io.spring.dependency-management' version '1.1.4'
}

group = 'com.petlog'
version = '1.0.0'
sourceCompatibility = '17'

configurations {
    compileOnly {
        extendsFrom annotationProcessor
    }
}

repositories {
    mavenCentral()
}

dependencies {
    // Spring Boot Starters
    implementation 'org.springframework.boot:spring-boot-starter-web'
    implementation 'org.springframework.boot:spring-boot-starter-data-jpa'
    implementation 'org.springframework.boot:spring-boot-starter-validation'
    implementation 'org.springframework.boot:spring-boot-starter-actuator'

    // Kafka
    implementation 'org.springframework.kafka:spring-kafka'

    // PostgreSQL
    runtimeOnly 'org.postgresql:postgresql'

    // AWS Bedrock SDK
    implementation platform('software.amazon.awssdk:bom:2.20.+')
    implementation 'software.amazon.awssdk:bedrock-runtime'
    implementation 'software.amazon.awssdk:bedrockruntime'
    implementation 'software.amazon.awssdk:auth'

    // Milvus Vector DB
    implementation 'io.milvus:milvus-sdk-java:2.3.4'

    // JSON Processing
    implementation 'com.fasterxml.jackson.core:jackson-databind'

    // Lombok
    compileOnly 'org.projectlombok:lombok'
    annotationProcessor 'org.projectlombok:lombok'
```

```
// Test
testImplementation 'org.springframework.boot:spring-boot-starter-test'
testImplementation 'org.springframework.kafka:spring-kafka-test'
testImplementation 'org.testcontainers:postgresql:1.19.3'
testImplementation 'org.testcontainers:junit-jupiter:1.19.3'
}

tasks.named('test') {
    useJUnitPlatform()
}
```

## 4.2 settings.gradle 작성

rootProject.name = 'healthcare-service'

---

## Step 5: Docker Compose로 인프라 구축

로컬 개발 환경에 필요한 인프라를 Docker Compose로 구성하겠습니다.

### docker-compose.yml 작성

프로젝트 루트에 생성:

version: '3.8'

```
services:
  # PostgreSQL (Healthcare DB)
  postgres:
    image: postgres:15-alpine
    container_name: healthcare-postgres
    environment:
      POSTGRES_DB: healthcare_db
      POSTGRES_USER: petlog
      POSTGRES_PASSWORD: petlog123
      POSTGRES_INITDB_ARGS: "--encoding=UTF-8"
    ports:
      - "5434:5432"
    volumes:
      - postgres-data:/var/lib/postgresql/data
    healthcheck:
      test: ["CMD-SHELL", "pg_isready -U petlog"]
      interval: 10s
      timeout: 5s
      retries: 5

  # Zookeeper (Kafka 의존성)
  zookeeper:
    image: confluentinc/cp-zookeeper:7.5.0
    container_name: healthcare-zookeeper
    environment:
      ZOOKEEPER_CLIENT_PORT: 2181
      ZOOKEEPER_TICK_TIME: 2000
    ports:
      - "2181:2181"

  # Kafka
  kafka:
    image: confluentinc/cp-kafka:7.5.0
    container_name: healthcare-kafka
    depends_on:
      - zookeeper
    ports:
      - "9092:9092"
    environment:
      KAFKA_BROKER_ID: 1
      KAFKA_ZOOKEEPER_CONNECT: zookeeper:2181
      KAFKA_ADVERTISED_LISTENERS: PLAINTEXT://localhost:9092
      KAFKA_OFFSETS_TOPIC_REPLICATION_FACTOR: 1
      KAFKA_TRANSACTION_STATE_LOG_MIN_ISR: 1
      KAFKA_TRANSACTION_STATE_LOG_REPLICATION_FACTOR: 1
    healthcheck:
      test: ["CMD", "kafka-broker-api-versions", "--bootstrap-server", "localhost:9092"]
      interval: 10s
      timeout: 10s
      retries: 5

  # Milvus - etcd
  etcd:
    image: quay.io/coreos/etcd:v3.5.5
    container_name: healthcare-milvus-etcd
    environment:
      - ETCD_AUTO_COMPACTION_MODE=revision
      - ETCD_AUTO_COMPACTION_RETENTION=1000
      - ETCD_QUOTA_BACKEND_BYTES=4294967296
      - ETCD_SNAPSHOT_COUNT=50000
    volumes:
      - etcd-data:/etc
    command: etcd --advertise-client-urls=http://127.0.0.1:2379 --listen-client-urls http://0.0.0.0:2379 --data-dir /etc
    healthcheck:
      test: ["CMD", "etcdctl", "endpoint", "health"]
      interval: 30s
      timeout: 20s
      retries: 3

  # Milvus - MinIO
  minio:
```

```

image: minio/minio:RELEASE.2023-03-20T20-16-18Z
container_name: healthcare-minio
environment:
  MINIO_ACCESS_KEY: minioadmin
  MINIO_SECRET_KEY: minioadmin
ports:
  - "9001:9000"
  - "9002:9001"
volumes:
  - minio-data:/minio_data
command: minio server /minio_data --console-address ":9001"
healthcheck:
  test: ["CMD", "curl", "-f", "http://localhost:9000/minio/health/live"]
  interval: 30s
  timeout: 20s
  retries: 3

# Milvus Standalone
milvus:
  image: milvusdb/milvus:v2.3.4
  container_name: healthcare-milvus
  depends_on:
    - etcd
    - minio
  ports:
    - "19530:19530"
    - "9091:9091"
  environment:
    ETCD_ENDPOINTS: etcd:2379
    MINIO_ADDRESS: minio:9000
  volumes:
    - milvus-data:/var/lib/milvus
  healthcheck:
    test: ["CMD", "curl", "-f", "http://localhost:9091/healthz"]
    interval: 30s
    start_period: 90s
    timeout: 20s
    retries: 3

volumes:
  postgres-data:
  etcd-data:
  minio-data:
  milvus-data:

```

## 인프라 실행

**Copy** # Docker Compose 실행 (백그라운드)  
docker-compose up -d

# 로그 확인  
docker-compose logs -f

# 서비스 상태 확인  
docker-compose ps

# 특정 서비스 로그만 보기  
docker-compose logs -f kafka

# 종료 (데이터 유지)  
docker-compose down

# 종료 (데이터 삭제)  
docker-compose down -v

## 실행 후 확인사항:

**Copy** # PostgreSQL 연결 테스트  
docker exec -it healthcare-postgres psql -U petlog -d healthcare\_db

# 연결되면:  
Wl # 데이터베이스 목록  
Wq # 종료

# Kafka 토픽 확인  
docker exec -it healthcare-kafka kafka-topics --list --bootstrap-server localhost:9092

## Step 6: Spring Boot Application 기본 설정

### 6.1 application.yml 작성

src/main/resources/application.yml:

**Copy** spring:
application:
 name: healthcare-service

profiles:
 active: local

jpa:
 hibernate:
 ddl-auto: update
 show-sql: true
 properties:
 hibernate:
 format\_sql: true

```
dialect: org.hibernate.dialect.PostgreSQLDialect
```

```
server:
  port: 8003
```

```
management:
  endpoints:
    web:
      exposure:
        include: health,info,metrics
  endpoint:
    health:
      show-details: always
```

## 6.2 application-local.yml 작성

src/main/resources/application-local.yml:

```
Copy spring:
  datasource:
    url: jdbc:postgresql://localhost:5434/healthcare_db
    username: petlog
    password: petlog123
    driver-class-name: org.postgresql.Driver

  kafka:
    bootstrap-servers: localhost:9092
    consumer:
      group-id: healthcare-group
      key-deserializer: org.apache.kafka.common.serialization.StringDeserializer
      value-deserializer: org.springframework.kafka.support.serializer.JsonDeserializer
      properties:
        spring.json.trusted.packages: "com.petlog.healthcare.*"
    producer:
      key-serializer: org.apache.kafka.common.serialization.StringSerializer
      value-serializer: org.springframework.kafka.support.serializer.JsonSerializer

  aws:
    bedrock:
      region: us-east-1
      # 로컬에서는 ~/.aws/credentials 사용

  milvus:
    host: localhost
    port: 19530

logging:
  level:
    com.petlog.healthcare: DEBUG
    org.springframework.kafka: INFO
```

---

## Step 7: Main Application 클래스 생성

src/main/java/com/petlog/healthcare/HealthcareServiceApplication.java:

```
Copy package com.petlog.healthcare;

import org.springframework.boot.SpringApplication;
import org.springframework.boot.autoconfigure.SpringBootApplication;
import org.springframework.kafka.annotation.EnableKafka;
import org.springframework.scheduling.annotation.EnableScheduling;

@SpringBootApplication
@EnableKafka
@EnableScheduling
public class HealthcareServiceApplication {

    public static void main(String[] args) {
        SpringApplication.run(HealthcareServiceApplication.class, args);
    }
}
```

---

## Step 8: 첫 번째 테스트 실행

### 8.1 Gradle Build

```
Copy # 프로젝트 루트에서
./gradlew clean build

# Windows라면:
gradlew.bat clean build
```

### 8.2 애플리케이션 실행

```
Copy # 방법 1: Gradle로 실행
./gradlew bootRun

# 방법 2: JAR 실행
java -jar build/libs/healthcare-service-1.0.0.jar

# 방법 3: IDE에서 실행
# HealthcareServiceApplication 클래스 우클릭 → Run
```

### 8.3 Health Check

브라우저 또는 curl로 확인:

```
Copy # Health Check
curl http://localhost:8003/actuator/health

# 기대 결과:
{
  "status": "UP",
  "components": {
    "db": {
      "status": "UP",
      "details": {
        "database": "PostgreSQL",
        "validationQuery": "isValid()"
      }
    },
    "diskSpace": {
      "status": "UP"
    },
    "ping": {
      "status": "UP"
    }
  }
}
```

---

## ✓ Phase 0 완료 체크리스트

- Java 17+ 설치 확인
  - AWS CLI 설치 및 자격 증명 설정
  - AWS Bedrock 모델 액세스 승인 (Claude, Titan)
  - Docker & Docker Compose 설치
  - 프로젝트 디렉토리 생성
  - build.gradle 작성
  - docker-compose.yml 작성 및 인프라 실행
  - application.yml 설정
  - Spring Boot 애플리케이션 실행 성공
  - Health Check 통과
- 

## ⌚ 다음 단계 안내

Phase 0이 완료되면 다음 순서로 진행합니다:

1. **Phase 1:** AWS Bedrock 연동 (Claude, Titan Embeddings)
  2. **Phase 2:** 건강 데이터 관리 (Entity, Repository, Service, Controller)
  3. **Phase 3:** Kafka Consumer 구현 (다이어리 이벤트 수신)
  4. **Phase 4:** Milvus Vector DB 연동 (벡터 저장/검색)
  5. **Phase 5:** 일반 챗봇 구현 (Claude Haiku)
  6. **Phase 6:** 페르소나 챗봇 구현 (Claude Sonnet + RAG)
- 

현재 상황을 알려주세요:

1. AWS Bedrock 액세스 키가 있으신가요?

- 있음 → AWS CLI 설정 진행
- 없음 → IAM 사용자 생성 방법 안내

2. 개발 환경은 어떤 OS를 사용하시나요?

- macOS / Windows / Linux

3. IDE는 무엇을 사용하시나요?

- IntelliJ IDEA / VSCode / Eclipse

4. 위 Step 중 막히는 부분이 있나요?

- 어느 부분에서 막히셨는지 알려주시면 상세히 설명드리겠습니다!

준비되시면 Phase 1부터 시작하겠습니다! 🚀

▣

복사

