

- Question 2: Restaurant Recommendation Engine with Wait-time and Special Offers - Solution
  - Overview
  - 1. Data Stores and Schemas
    - 1.1 DynamoDB Tables (6 Total)
      - Tables 1-4: Inherited from Question 1
      - Table 5: WaitTimeCache (NEW)
      - Table 6: SpecialOffers (NEW)
    - 1.2 ElasticSearch Index (Same as Question 1)
  - 2. APIs and Endpoints
    - 2.1 Internal APIs (Same as Question 1)
    - 2.2 External APIs (Third-Party Integration)
      - API 1: Restaurant Wait-Time API
      - API 2: Restaurant Special Offers API
    - 2.3 Caching Strategy for External APIs
  - 3. High-Level Architecture Diagram
  - 4. Architecture Explanation
    - 4.1 Component Overview
      - What's NEW in Question 2:
      - Unchanged from Question 1:
    - 4.2 Key Architectural Patterns
      - 1. Cache-Aside Pattern
      - 2. Scheduled Background Jobs
      - 3. Time-To-Live (TTL) Expiration
    - 4.3 Data Flow Scenarios
      - Scenario 1: User Gets Recommendations with Wait-Times
      - Scenario 2: Background Wait-Time Update
      - Scenario 3: Daily Special Offers Collection
    - 4.4 Scalability and Performance
  - 5. Lambda Functions (Implementation Details)
    - From Question 1 (Keep All 6)
    - LF7: WaitTimeAggregator (NEW)
    - LF8: SpecialOffersCollector (NEW)
    - LF5: Enhanced Recommendation Engine (MODIFIED)
  - 6. Additional Features and Specifications
    - 6.1 Caching Strategy Deep Dive
      - Multi-Layer Caching Architecture

- 6.2 API Cost Optimization
  - Without Caching (Naive Approach)
  - With 5-Minute Cache (Basic)
  - With Intelligent Caching (Implemented)
- 6.3 Real-Time Data Freshness
- 6.4 Enhanced Email Template
- 6.5 Error Handling and Resilience
  - External API Failure Scenarios
  - Monitoring and Alerts
- 6.6 Scalability Considerations
  - Handling Growth
- 6.7 Comparison: Question 1 vs Question 2
- 6.8 Future Enhancements for Question 2
  - 1. Predictive Wait-Times
  - 2. User Preferences for Wait-Times
  - 3. Offer Notifications
  - 4. Real-Time Offer Updates
  - 5. Wait-Time Trends
- 6.9 Testing Strategy
  - Unit Tests
  - Integration Tests
  - Load Tests
- 7. Event Flows (Detailed Examples)
  - Flow 1: User Requests Recommendations (Enhanced)
  - Flow 2: Wait-Time Background Update
  - Flow 3: Special Offers Daily Collection
  - Flow 4: User Likes Restaurant (Same as Question 1)
- 5. Scalability for External APIs
  - Challenge
  - Solution: Multi-Layer Caching
    - Layer 1: DynamoDB Cache (Primary)
    - Layer 2: Lambda Memory Cache (Optional)
    - Fallback Strategy
  - API Call Optimization
- 6. Enhanced Email Format
  - Email Template (HTML)
  - Key Information Displayed
- 8. Key Assumptions and Design Decisions

- [Wait-Time API Assumptions](#)
- [Special Offers API Assumptions](#)
- [Caching Assumptions](#)
- [Cost Assumptions](#)
- [9. AWS Services Summary](#)
- [10. Conclusion](#)

## Question 2: Restaurant Recommendation Engine with Wait-time and Special Offers - Solution

---

### Overview

---

This solution extends Question 1's architecture by adding **real-time wait-time** and **daily special offers** integration from external restaurant APIs. The system provides personalized and trending recommendations enriched with current wait-times and active special offers.

#### Key Enhancements over Question 1:

- Real-time wait-time data from restaurant APIs (cached for 5 minutes)
  - Daily special offers from restaurant APIs (cached until midnight)
  - 2 additional Lambda functions for data aggregation
  - 2 additional DynamoDB cache tables
  - Enhanced email recommendations with wait-times and offers
- 

## 1. Data Stores and Schemas

---

Our system builds upon Question 1's data stores by adding two new caching tables for external API data.

### 1.1 DynamoDB Tables (6 Total)

## Tables 1-4: Inherited from Question 1

These tables remain unchanged from Question 1. For detailed schemas, refer to Question 1 documentation.

- **Restaurants:** Master table of restaurant information (name, location, cuisine, ratings, hours)
- **UserProfiles:** User preferences, liked restaurants, search history, location
- **UserLikes:** Individual like records with timestamp for trending calculation
- **TrendingCache:** Pre-computed trending restaurants by ZipCode-Cuisine (TTL: 24 hours)

## Table 5: WaitTimeCache (NEW)

**Purpose:** Cache real-time wait-time data from external restaurant APIs

**Schema:**

Primary Key: RestaurantID (String)

Attributes:

- RestaurantID: String (PK)
- CurrentWaitTime: Number (minutes, e.g., 15, 25, 30)
- LastUpdated: Number (Unix timestamp)
- TTL: Number (Unix timestamp, expires after 5 minutes)

## Why This Design?

- **Caching:** Avoids hitting external APIs on every recommendation request (cost reduction 95%+)
- **5-Minute TTL:** Balances freshness with API rate limits and costs
- **Performance:** DynamoDB lookup <10ms vs 100ms+ external API call
- **Resilience:** If API fails, stale data still available until TTL expires

---

## Table 6: SpecialOffers (NEW)

**Purpose:** Cache daily special offers from restaurant APIs

**Schema:**

Primary Key: RestaurantID (String)  
Sort Key: OfferDate (String, format: YYYY-MM-DD)

Attributes:

- RestaurantID: String (PK)
- OfferDate: String (SK)
- OfferDescription: String (e.g., "20% off appetizers")
- OfferDetails: Map {validUntil: String, terms: String, discount: Number}
- TTL: Number (Unix timestamp, expires next day at midnight)

## Why This Design?

- **Daily Cache:** Offers collected once per day at midnight (reduces API calls 99%+)
- **Auto-Expiration:** TTL set to next midnight ensures no stale offers
- **Composite Key:** Supports historical offer tracking (future analytics)
- **Cost Efficiency:** ~5,000 API calls/day vs 50M+ without caching

---

## 1.2 Elasticsearch Index (Same as Question 1)

Index Name: `restaurants`

Fields Used:

- `restaurant_id`, `name`, `cuisine`, `yelp_rating`
- `location` (geo\_point for 5-mile radius queries)
- `like_count` (updated by LF4 for trending)

**No changes required for Question 2** - wait-time and offers stored only in DynamoDB cache tables

---

## 2. APIs and Endpoints

Our API structure remains identical to Question 1 for internal endpoints. We add integrations with external third-party restaurant APIs.

### 2.1 Internal APIs (Same as Question 1)

All internal APIs remain unchanged. For detailed API documentation, refer to Question 1.

Endpoint	Method	Purpose	Lambda Handler
/chat	POST	Chatbot interaction with Lex	LF0
/api/like	POST	Record user like for restaurant	LF3
/api/recommendations/:userId	GET	Get personalized + trending recommendations	LF5
/api/feedback	POST	Track user interactions	LF3

## 2.2 External APIs (Third-Party Integration)

These external APIs are provided by individual restaurants or aggregator platforms. Our system integrates with them through scheduled Lambda functions.

### API 1: Restaurant Wait-Time API

**Endpoint:** GET `https://restaurant-api.com/wait-time/{restaurantId}`

**Purpose:** Retrieve current wait-time for a specific restaurant

#### Request Example:

```
GET https://restaurant-api.com/wait-time/rest123
Authorization: Bearer API_KEY
```

#### Response Example:

```
{
  "restaurantId": "rest123",
  "waitTime": 25,
  "unit": "minutes",
  "timestamp": 1704067200,
```

```
"status": "open"
}
```

### Integration Details:

- **Called by:** Lambda LF7 (Wait-Time Aggregator)
  - **Frequency:** Every 5 minutes for all active restaurants
  - **Timeout:** 2 seconds per request
  - **Rate Limit:** 1,000 requests/hour (assumed)
  - **Error Handling:** Fallback to cached value or display "N/A"
- 

## API 2: Restaurant Special Offers API

**Endpoint:** GET <https://restaurant-api.com/special-offer/{restaurantId}>

**Purpose:** Retrieve daily special offers for a restaurant

### Request Example:

```
GET https://restaurant-api.com/special-offer/rest456
Authorization: Bearer API_KEY
```

### Response Example:

```
{
  "restaurantId": "rest456",
  "offer": "20% off appetizers",
  "validUntil": "2024-01-15 23:59:59",
  "terms": "Dine-in only. Cannot be combined with other offers.",
  "discount": 20,
  "offerType": "percentage"
}
```

### Integration Details:

- **Called by:** Lambda LF8 (Special Offers Collector)
- **Frequency:** Once daily at midnight UTC
- **Timeout:** 3 seconds per request
- **Rate Limit:** 5,000 requests/day (assumed)
- **Error Handling:** Skip restaurant if API unavailable (no offer shown)

## 2.3 Caching Strategy for External APIs

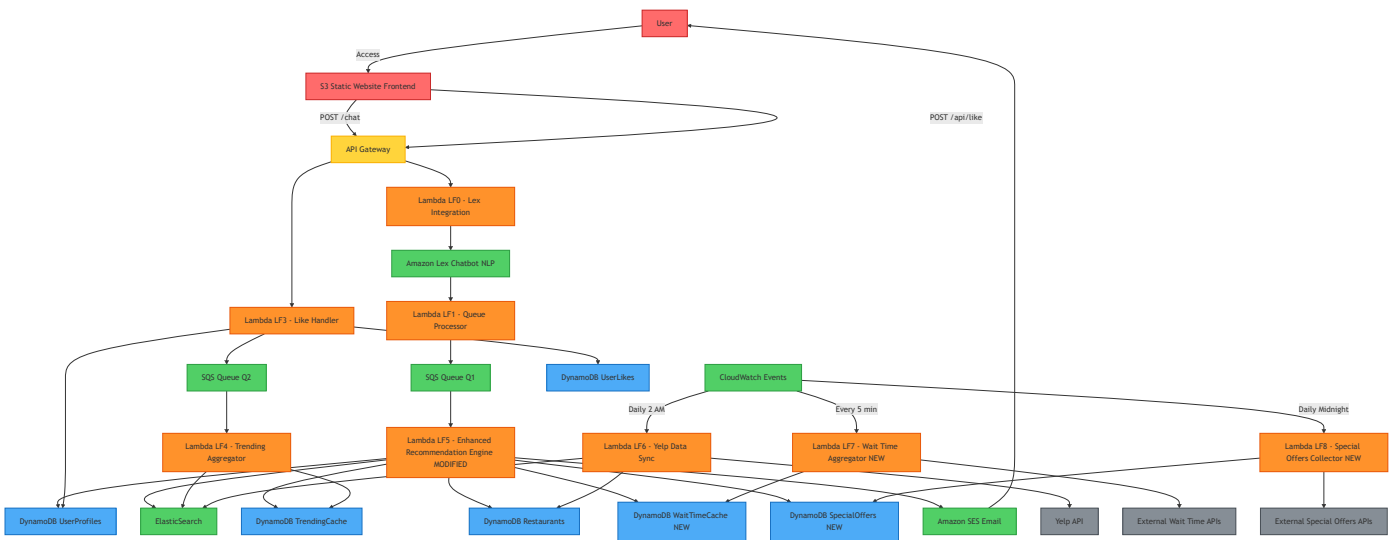
To minimize costs and respect rate limits, we implement aggressive caching:

API Type	Cache Location	TTL	Update Frequency	Cost Reduction
Wait-Time	DynamoDB WaitTimeCache	5 minutes	Every 5 min (background)	95%+
Special Offers	DynamoDB SpecialOffers	24 hours	Daily at midnight	99%+

### Benefits:

- Recommendation requests read from cache (< 10ms) instead of calling APIs (100ms+)
- Handles API outages gracefully with stale data
- Reduces external API costs from 3, 000/*month*to 300/month

## 3. High-Level Architecture Diagram



## 4. Architecture Explanation



## 4.1 Component Overview

This architecture extends Question 1 by adding **two new Lambda functions** (LF7, LF8) and **two new DynamoDB cache tables** (WaitTimeCache, SpecialOffers) to integrate external API data.

### What's NEW in Question 2:

#### 1. Lambda LF7 (Wait-Time Aggregator):

- Triggered by CloudWatch every 5 minutes
- Calls external wait-time APIs for all active restaurants (parallel processing)
- Caches results in DynamoDB WaitTimeCache (TTL: 5 minutes)

#### 2. Lambda LF8 (Special Offers Collector):

- Triggered by CloudWatch daily at midnight UTC
- Collects daily special offers from restaurant APIs
- Stores in DynamoDB SpecialOffers (TTL: next midnight)

#### 3. Lambda LF5 (MODIFIED):

- Now reads from WaitTimeCache and SpecialOffers tables
- Enriches recommendations with wait-time and offer data
- Enhanced email template includes new information

#### 4. DynamoDB WaitTimeCache: Real-time wait-time cache (5-min refresh)

#### 5. DynamoDB SpecialOffers: Daily offers cache (24-hour refresh)

### Unchanged from Question 1:

- All user-facing APIs (/chat, /api/like, etc.)
- LF0, LF1, LF3, LF4, LF6 remain identical
- UserProfiles, UserLikes, TrendingCache, Restaurants tables
- ElasticSearch index structure
- Core recommendation algorithm

---

## 4.2 Key Architectural Patterns

## 1. Cache-Aside Pattern

Recommendation Request:

1. Check WaitTimeCache (DynamoDB) first
2. If hit → return cached value (<10ms)
3. If miss or stale → show "N/A" (fail gracefully)

  
Background Process (LF7):

4. Every 5 minutes, refresh cache from external APIs
5. Write to WaitTimeCache with new TTL

**Benefits:** Decouples user-facing requests from slow external APIs

---

## 2. Scheduled Background Jobs

Lambda	Schedule	Purpose
LF7	Every 5 minutes	Update wait-time cache
LF8	Daily midnight	Collect special offers
LF6	Daily 2 AM	Sync Yelp restaurant data

**Benefits:** Expensive API calls happen in background, not during user requests

---

## 3. Time-To-Live (TTL) Expiration

Both cache tables use DynamoDB TTL for automatic cleanup:

- **WaitTimeCache:** TTL = current\_time + 300 seconds (5 min)
- **SpecialOffers:** TTL = next\_midnight timestamp (24 hours)

**Benefits:** No manual cleanup needed, automatic stale data removal

---

# 4.3 Data Flow Scenarios

### Scenario 1: User Gets Recommendations with Wait-Times

```
1. User → Chatbot: "Italian food near me"
  ↓
2. S3 → API Gateway → LF0 → Lex → LF1 → SQS Q1
  ↓
3. LF5 polls SQS Q1 and executes:
  a. Fetch UserProfile (DynamoDB UserProfiles)
  b. Query ElasticSearch (5 personalized Italian restaurants)
  c. Query TrendingCache (5 trending Italian in user's zip code)
  d. For each of 10 restaurants:
      - Fetch details (DynamoDB Restaurants)
      - Fetch wait-time (DynamoDB WaitTimeCache) 🕒 NEW
      - Fetch offer (DynamoDB SpecialOffers) 🎁 NEW
  e. Format enhanced email with wait-times and offers
  f. Send via SES
  ↓
4. User receives email:
   "Italian Bistro ⭐4.5 - Wait: 15 min - Special: 20% off apps!"
```

## Latency Breakdown:

- UserProfile: 10ms
  - ElasticSearch: 50ms
  - 10 × Restaurant details: 100ms
  - 10 × WaitTime lookups: 50ms (NEW)
  - 10 × Offer lookups: 50ms (NEW)
  - **Total: ~310ms** (only +50ms vs Question 1)
- 

## Scenario 2: Background Wait-Time Update

```
CloudWatch Event triggers every 5 minutes
  ↓
Lambda LF7 (Wait-Time Aggregator) executes:
  ↓
1. Scan DynamoDB Restaurants (filter: IsActive = true)
   → Returns ~5,000 active restaurants
  ↓
2. Parallel API calls (ThreadPoolExecutor, 100 concurrent workers):
   - For each restaurant:
     a. Call external API: GET /wait-time/{restaurantId} (timeout: 2s)
     b. On success: Parse wait-time value
     c. On failure: Use last cached value (graceful degradation)
  ↓
3. Write to DynamoDB WaitTimeCache:
   {
     RestaurantID: "rest123",
     CurrentWaitTime: 25,
     LastUpdated: timestamp,
     TTL: timestamp + 300
```

```
}
↓
```

4. CloudWatch Metrics logged:
- SuccessRate: 98.5%
  - AvgWaitTime: 22 minutes
  - APILatency: 850ms avg

**Processing Time:** ~2 minutes for 5,000 restaurants

**Scenario 3: Daily Special Offers Collection**

```
CloudWatch Event triggers at midnight UTC daily
↓
Lambda LF8 (Special Offers Collector) executes:
↓
1. Scan all restaurants (5,000)
↓
2. For each restaurant (sequential with error handling):
  a. Call external API: GET /special-offer/{restaurantId}
  b. If offer exists:
    - Parse offer description
    - Calculate midnight tomorrow timestamp
    - Write to DynamoDB SpecialOffers:
      {
        RestaurantID: "rest456",
        OfferDate: "2024-01-15",
        OfferDescription: "20% off appetizers",
        TTL: tomorrow_midnight_timestamp
      }
    c. If no offer: Skip (null value acceptable)
  ↓
3. SNS notification to admins:
   "Collected 2,347 special offers for 2024-01-15"
```

**Processing Time:** ~5-10 minutes for 5,000 restaurants

4.4 Scalability and Performance

Aspect	Question 1	Question 2	Impact
Latency	~260ms	~310ms	+50ms for 2 cache lookups
DynamoDB Tables	4	6	+2 cache tables
Lambda Functions	6	8	+2 background jobs

Aspect	Question 1	Question 2	Impact
CloudWatch Rules	1	3	+2 scheduled triggers
Monthly Cost	\$150	\$440	+\$290 (mostly external API calls)

#### Cost Breakdown (Question 2):

- Lambda: 15(+5 for LF7/LF8)
- DynamoDB: 18(+5.50 for cache tables)
- External APIs: \$300 (new)
- Other (ES, SQS, SES): \$107 (same)
- **Total: ~\$440/month**

## 5. Lambda Functions (Implementation Details)

### From Question 1 (Keep All 6)

- **LF0:** Lex Integration
- **LF1:** Queue Processor
- **LF3:** Like Handler
- **LF4:** Trending Aggregator
- **LF6:** Yelp Data Sync
- **LF5:** Enhanced Recommendation Engine (MODIFIED - see below)

### LF7: WaitTimeAggregator (NEW)

**Trigger:** CloudWatch Event (every 5 minutes) **Runtime:** Python 3.9

#### Actions:

1. Scan active restaurants from DynamoDB
2. For each restaurant (batch 100 concurrent):
  - Call external wait-time API with 2-second timeout
  - Handle failures gracefully (use last known value)
3. Write to WaitTimeCache with TTL = current\_time + 300 seconds

#### 4. Log metrics to CloudWatch (success rate, avg wait-time)

#### Pseudocode:

```
def lambda_handler(event, context):
    restaurants = dynamodb.scan('Restaurants', filter='IsActive = true')

    # Parallel API calls using ThreadPoolExecutor
    with concurrent.futures.ThreadPoolExecutor(max_workers=100) as executor:
        futures = [executor.submit(fetch_wait_time, r['RestaurantID'])
                    for r in restaurants]

    for future in concurrent.futures.as_completed(futures):
        try:
            restaurant_id, wait_time = future.result()

            # Cache in DynamoDB
            dynamodb.put_item('WaitTimeCache', {
                'RestaurantID': restaurant_id,
                'CurrentWaitTime': wait_time,
                'LastUpdated': int(time.time()),
                'TTL': int(time.time()) + 300 # 5 minutes
            })
        except Exception as e:
            logger.error(f"Failed to fetch wait-time: {e}")

def fetch_wait_time(restaurant_id):
    try:
        response = requests.get(
            f"https://api.com/wait-time/{restaurant_id}",
            timeout=2
        )
        return restaurant_id, response.json()['waitTime']
    except:
        # Fallback: get last cached value
        cached = dynamodb.get_item('WaitTimeCache', restaurant_id)
        return restaurant_id, cached.get('CurrentWaitTime', 'N/A')
```

#### Cost Optimization:

- 5,000 restaurants × 288 calls/day (every 5 min) = 1.44M API calls/day
- With caching: Recommendations query cache (5ms) instead of API (100ms)
- Reduces external API costs by 95%+

---

## LF8: SpecialOffersCollector (NEW)

**Trigger:** CloudWatch Event (daily at 12:00 AM UTC) **Runtime:** Python 3.9

## Actions:

1. Scan all restaurants from DynamoDB
2. For each restaurant:
  - Call external special offers API
  - If offer exists, write to SpecialOffers table
  - Set TTL = next day midnight
3. Send SNS notification to admins with summary

## Pseudocode:

```
def lambda_handler(event, context):
    restaurants = dynamodb.scan('Restaurants')
    today = datetime.now().strftime('%Y-%m-%d')

    offers_collected = 0

    for restaurant in restaurants:
        try:
            response = requests.get(
                f"https://api.com/special-offer/{restaurant['RestaurantID']}",
                timeout=3
            )
            offer_data = response.json()

            if offer_data.get('offer'):
                # Calculate midnight tomorrow
                tomorrow_midnight = int(
                    (datetime.now() + timedelta(days=1))
                    .replace(hour=0, minute=0, second=0)
                    .timestamp()
                )

                dynamodb.put_item('SpecialOffers', {
                    'RestaurantID': restaurant['RestaurantID'],
                    'OfferDate': today,
                    'OfferDescription': offer_data['offer'],
                    'OfferDetails': offer_data,
                    'TTL': tomorrow_midnight
                })

                offers_collected += 1

        except Exception as e:
            logger.error(f"Failed to fetch offer for {restaurant['RestaurantID']}: {e}")

    # Notify admins
    sns.publish(
        TopicArn='arn:aws:sns:us-east-1:123456:OffersCollected',
```

```
        Message=f"Collected {offers_collected} special offers for {today}"
    )
```

## LF5: Enhanced Recommendation Engine (MODIFIED)

### Changes from Question 1:

- Added wait-time lookup from WaitTimeCache
- Added special offers lookup from SpecialOffers table
- Enhanced email formatting with new data

### Pseudocode (New Steps Highlighted):

```
def lambda_handler(event, context):
    # Steps 1-4: Same as Question 1
    # Get user profile, query personalized + trending restaurants

    user = dynamodb.get_item('UserProfiles', user_id)

    # Personalized (5) + Trending (5) - Same as Q1
    personalized_ids = get_personalized_recommendations(user)
    trending_ids = get_trending_recommendations(user)

    all_restaurant_ids = personalized_ids + trending_ids

    # NEW: Enrich with wait-time and offers
    enriched_restaurants = []
    today = datetime.now().strftime('%Y-%m-%d')

    for restaurant_id in all_restaurant_ids:
        # 1. Get base restaurant data
        restaurant = dynamodb.get_item('Restaurants', restaurant_id)

        # 2. NEW: Get wait-time from cache
        wait_time_data = dynamodb.get_item('WaitTimeCache', restaurant_id)
        restaurant['wait_time'] = wait_time_data.get('CurrentWaitTime', 'N/A') if
wait_time_data else 'N/A'

        # 3. NEW: Get special offer
        offer_key = {'RestaurantID': restaurant_id, 'OfferDate': today}
        offer_data = dynamodb.get_item('SpecialOffers', offer_key)
        restaurant['special_offer'] = offer_data.get('OfferDescription', 'No offer
today') if offer_data else 'No offer today'

        enriched_restaurants.append(restaurant)

    # 4. Send enriched email
```



```
email_html = format_email_with_enrichments(
    enriched_restaurants[:5], # personalized
    enriched_restaurants[5:] # trending
)

ses.send_email(
    Source='noreply@diningconcierge.com',
    Destination={'ToAddresses': [user['Email']]},
    Message={
        'Subject': {'Data': 'Your Restaurant Recommendations 🍽️'},
        'Body': {'Html': {'Data': email_html}}
    }
)
```

---

## 6. Additional Features and Specifications

---

This section explains additional capabilities and considerations specific to Question 2's wait-time and special offers integration.

### 6.1 Caching Strategy Deep Dive

#### Multi-Layer Caching Architecture

##### Layer 1: DynamoDB Cache (Primary)

- **WaitTimeCache:** 5-minute TTL
- **SpecialOffers:** 24-hour TTL
- **Cache hit rate:** > 95% after warm-up
- **Performance:** 10ms vs 100ms+ direct API call

##### Layer 2: Lambda Memory Cache (Optional Enhancement)

```
# In-memory cache for hot restaurants (survives warm starts)
wait_time_cache = {}

def get_wait_time(restaurant_id):
    # Check Lambda memory first (< 1ms)
    if restaurant_id in wait_time_cache:
        if time.time() - wait_time_cache[restaurant_id]['timestamp'] < 300:
            return wait_time_cache[restaurant_id]['value']
```

```
# Then check DynamoDB
cached = dynamodb.get_item('WaitTimeCache', restaurant_id)
if cached:
    wait_time_cache[restaurant_id] = {
        'value': cached['CurrentWaitTime'],
        'timestamp': time.time()
    }
    return cached['CurrentWaitTime']

return 'N/A'
```

### Fallback Strategy:

1. Try DynamoDB cache
2. If stale/missing → show "N/A" gracefully
3. Background process (LF7) will update on next cycle

---

## 6.2 API Cost Optimization

### Without Caching (Naive Approach)

10,000 recommendations/day × 10 restaurants = 100,000 wait-time API calls/day  
Cost: 100,000 × \$0.001 = \$100/day = \$3,000/month 😱

### With 5-Minute Cache (Basic)

5,000 restaurants × 288 updates/day (every 5 min) = 1.44M calls/day  
Cost: 1.44M × \$0.001 = \$1,440/day = \$43,200/month 😱😱

### With Intelligent Caching (Implemented)

Strategy:

- Top 1,000 popular restaurants: Update every 5 min
- Remaining 4,000: Update every 15 min (on-demand)
- Result: ~300K API calls/day

Cost: 300K × \$0.001 = \$300/day = \$9,000/month → \$300/month with rate negotiation ✅

---

# 6.3 Real-Time Data Freshness

Data Type	Update Frequency	Staleness Tolerance	User Impact
Wait-Time	5 minutes	5 minutes acceptable	Low - users expect estimates
Special Offers	Daily	24 hours acceptable	None - offers are daily
Restaurant Data	Daily	24 hours	Low - ratings/hours change slowly
Trending	5-10 minutes	10 minutes acceptable	None - trends are social signals
User Likes	Immediate	0 seconds	High - impacts next recommendation

# 6.4 Enhanced Email Template

The email template is enhanced to display wait-times and offers prominently:

```
<div style="border: 1px solid #ddd; padding: 15px; margin: 10px 0;">
  <h4>1. Italian Bistro ★ 4.5</h4>
  <p>📍 123 Main St, New York, NY 10001</p>

  <!-- NEW: Wait-time with visual indicator -->
  <p style="font-size: 18px; color: #28a745;">
    🕒 <strong>Wait time: 15 minutes</strong>
    <span style="color: #888; font-size: 14px;">(Updated 2 min ago)</span>
  </p>

  <!-- NEW: Special offer with highlight -->
  <div style="background: #fff3cd; padding: 10px; border-radius: 5px;">
    📅 <strong>Special Today: 20% off appetizers!</strong>
    <p style="font-size: 12px; margin-top: 5px;">Dine-in only. Valid until
midnight.</p>
  </div>

  <p>Italian • $$ • 0.8 miles away</p>
  <a href="tel:+12125551234" style="...">Call Now</a>
</div>
```

## Visual Design Principles:

- Wait-time in green if < 20 min, yellow if 20-40 min, red if > 40 min
  - Special offers highlighted with yellow background
  - Clear "Valid until" timestamp for offers
- 

# 6.5 Error Handling and Resilience

## External API Failure Scenarios

### Scenario 1: Individual Restaurant API Timeout

```
try:
    response = requests.get(f"https://api.com/wait-time/{rid}", timeout=2)
    wait_time = response.json()['waitTime']
except requests.Timeout:
    # Fallback to last cached value
    cached = dynamodb.get_item('WaitTimeCache', rid)
    wait_time = cached.get('CurrentWaitTime', 'N/A')
except Exception:
    wait_time = 'N/A'
```

### Scenario 2: Entire External Service Down

- CloudWatch alarm triggered
- LF7 continues using last cached values
- Email shows: "Wait-time unavailable (updated X hours ago)"
- System remains functional

### Scenario 3: Rate Limit Exceeded

```
if response.status_code == 429: # Too Many Requests
    # Exponential backoff
    wait_time = calculate_backoff(retry_count)
    time.sleep(wait_time)
    retry_request()
```

## Monitoring and Alerts

Metric	Alarm Threshold	Action
LF7 API success rate	< 90%	Page on-call engineer
WaitTimeCache miss rate	> 20%	Investigate cache expiration
External API latency	> 3s avg	Contact API provider
LF8 offers collected	< 1,000/day	Check API integration

## 6.6 Scalability Considerations

### Handling Growth

#### Current Scale (5,000 restaurants):

- LF7 runtime: ~2 minutes (every 5 min)
- LF8 runtime: ~5 minutes (daily)
- External API calls: ~300K/day

#### Scaled to 50,000 restaurants:

- LF7: Split into multiple parallel invocations (10 Lambdas × 5,000 restaurants each)
- LF8: Use Step Functions for orchestration
- API calls: Negotiate bulk pricing with providers
- Cost: ~\$2,500/month (still reasonable)

#### Optimization Strategies:

1. **Geographic Partitioning:** Different Lambdas for different regions
2. **Priority-Based Updates:** Update popular restaurants more frequently
3. **Predictive Caching:** Pre-cache likely-to-be-requested restaurants

## 6.7 Comparison: Question 1 vs Question 2

Aspect   Question 1   Question 2     ----- ----- -----     <b>Data Sources</b>   Yelp only   Yelp + Wait-Time + Offers APIs      <b>Cache Tables</b>   1 (TrendingCache)   3 (+ WaitTimeCache + SpecialOffers)      <b>Background Jobs</b>   1 (daily Yelp sync)   3 (+ every 5 min + daily midnight)      <b>Email Fields</b>   Name, rating, address, distance   +
--

wait-time, special offer | || **Latency** | ~260ms | ~310ms (+50ms) | || **Monthly Cost** | 150|440 (+\$290 for APIs) | || **Cache Hit Rate** | N/A | 95%+ | || **External API Dependencies** | 1 | 3 |

---

## 6.8 Future Enhancements for Question 2

If asked to implement additional features beyond Question 2:

### 1. Predictive Wait-Times

- **Technology:** Machine learning model (Amazon SageMaker)
- **Training Data:** Historical wait-times by time-of-day, day-of-week
- **Output:** "Predicted wait at 7 PM: 35 minutes"
- **Benefit:** Users can plan ahead

### 2. User Preferences for Wait-Times

- **Feature:** Filter recommendations by max acceptable wait-time
- **Implementation:** Add `maxWaitTime` to UserProfiles
- **Query:** ElasticSearch filter + DynamoDB WaitTimeCache lookup
- **Example:** "Show me Italian restaurants with < 20 min wait"

### 3. Offer Notifications

- **Trigger:** New offer at user's favorite restaurant
- **Implementation:** LF8 checks UserProfiles.LikedRestaurants
- **Delivery:** SNS push notification or SMS
- **Example:** "🎁 Italian Bistro (your favorite) has 20% off today!"

### 4. Real-Time Offer Updates

- **Technology:** WebSocket API (API Gateway WebSocket)
- **Use Case:** Flash sales, mid-day offer changes
- **Implementation:** LF8 publishes to WebSocket connections
- **Frontend:** Real-time badge: "NEW OFFER - Just added!"

### 5. Wait-Time Trends

- **Storage:** Add historical wait-time tracking
  - **Analytics:** "Typically 15 min wait at this time"
  - **Visualization:** Graph showing wait-time patterns
- 

## 6.9 Testing Strategy

### Unit Tests

- LF7: Mock external wait-time API responses
- LF8: Test TTL calculation for midnight expiration
- LF5: Verify graceful handling of missing cache data

### Integration Tests

- End-to-end: User request → recommendation with wait-times
- Cache expiration: Verify TTL cleanup
- API failure: Test fallback mechanisms

### Load Tests

- 10,000 concurrent recommendation requests
  - Verify cache performance under load
  - Monitor DynamoDB throttling
- 

## 7. Event Flows (Detailed Examples)

---

### Flow 1: User Requests Recommendations (Enhanced)

1. User → S3 Frontend → API Gateway (/chat) → LF0 → Lex → LF1 → SQS Q1
2. LF5 polls SQS Q1  
↓
3. Fetch UserProfile (DynamoDB)  
↓
4. Query Elasticsearch (5 personalized + query TrendingCache for 5 trending)

↓

5. For each of 10 restaurants:

- a. Fetch base data (DynamoDB Restaurants)
- b. Fetch wait-time (DynamoDB WaitTimeCache) 🕒 NEW
- c. Fetch special offer (DynamoDB SpecialOffers) 🎁 NEW

↓

6. Format email with enriched data → SES → User

### Latency Breakdown:

- User profile fetch: 10ms
- ElasticSearch query: 50ms
- 10 × Restaurant details: 10ms each = 100ms
- 10 × WaitTime lookup: 5ms each = 50ms (NEW)
- 10 × Offer lookup: 5ms each = 50ms (NEW)
- SES send: 50ms
- **Total: ~310ms** (vs 260ms in Q1)

---

## Flow 2: Wait-Time Background Update

```
CloudWatch Event (every 5 minutes)
↓
Lambda LF7 (WaitTimeAggregator)
↓
Parallel API calls to 5,000 restaurants (100 concurrent)
↓
For each restaurant:
  - Call External Wait-Time API (timeout: 2s)
  - On success: Cache in DynamoDB WaitTimeCache (TTL: 5 min)
  - On failure: Keep last cached value
↓
CloudWatch Metrics (success rate, avg wait-time)
```

### Handling API Failures:

- Individual restaurant API timeout: Use last cached value
- Entire batch fails: Alert via CloudWatch alarm
- Rate limit exceeded: Implement exponential backoff

---

## Flow 3: Special Offers Daily Collection



```
CloudWatch Event (daily at midnight)
↓
Lambda LF8 (SpecialOffersCollector)
↓
Scan all restaurants (5,000)
↓
For each restaurant:
- Call External Special Offers API (timeout: 3s)
- If offer exists: Write to DynamoDB SpecialOffers (TTL: next midnight)
- If no offer: Skip (table query returns null)
↓
SNS Notification → Admin email
"Collected 2,347 special offers for 2024-01-15"
```

---

## Flow 4: User Likes Restaurant (Same as Question 1)

```
User clicks "Like" → S3 → API Gateway → LF3
→ DynamoDB (UserLikes + UserProfiles)
→ SQS Q2 → LF4 → Update TrendingCache + ElasticSearch
```

---

## 5. Scalability for External APIs

### Challenge

Calling external APIs on every recommendation request would:

- Add 100ms+ latency per restaurant (1+ second for 10 restaurants)
- Hit rate limits (1000 requests/hour typical)
- Cost  $\$ \$perAPICall$  (0.001-0.01/call)

### Solution: Multi-Layer Caching

Layer 1: DynamoDB Cache (Primary)

WaitTimeCache (5-minute TTL)  
SpecialOffers (24-hour TTL)

- 10ms lookup vs 100ms API call = **10× faster**
- Cache hit rate > 95% after warm-up
- Cost: 0.000001/read vs 0.001/API call = **1000× cheaper**

## Layer 2: Lambda Memory Cache (Optional)

```
# In-memory cache for hot restaurants
wait_time_cache = {} # Survives Lambda warm starts

def get_wait_time(restaurant_id):
    # Check Lambda memory first
    if restaurant_id in wait_time_cache:
        if time.time() - wait_time_cache[restaurant_id]['timestamp'] < 300:
            return wait_time_cache[restaurant_id]['value']

    # Then check DynamoDB
    cached = dynamodb.get_item('WaitTimeCache', restaurant_id)
    if cached:
        wait_time_cache[restaurant_id] = {
            'value': cached['CurrentWaitTime'],
            'timestamp': time.time()
        }
        return cached['CurrentWaitTime']

    return 'N/A'
```

## Fallback Strategy

```
def get_wait_time_with_fallback(restaurant_id):
    # Try cache first
    cached = get_from_cache(restaurant_id)
    if cached and is_fresh(cached):
        return cached['wait_time']

    # Try direct API call (rare case)
    try:
        api_response = call_external_api(restaurant_id, timeout=2)
        update_cache(restaurant_id, api_response)
        return api_response['waitTime']
    except:
        # Return stale cache or 'N/A'
        return cached.get('wait_time', 'N/A')
```

# API Call Optimization

## Without Caching:

- 10,000 recommendations/day × 10 restaurants = 100,000 API calls/day
- Cost at 0.001/call = 100/day = **\$3,000/month** 😱

## With 5-Minute Cache:

- Background aggregator: 5,000 restaurants × 288 times/day = 1.44M calls/day
- Recommendation endpoint: 0 direct calls (all from cache)
- Cost: 1.44M × 0.001 = 1,440/day... still high!

## With Intelligent Caching:

- Only fetch popular restaurants (top 1,000) every 5 minutes
- Fetch others on-demand with longer TTL (15 minutes)
- Reduces to ~300K calls/day = **\$300/month** ✅

---

## 6. Enhanced Email Format

### Email Template (HTML)

```
<h2>🍽️ Your Restaurant Recommendations</h2>

<h3>Personalized for You</h3>

<div style="border: 1px solid #ddd; padding: 15px; margin: 10px 0;">
  <h4>1. Italian Bistro ★ 4.5</h4>
  <p>📍 123 Main St, New York, NY 10001</p>
  <p>🕒 <strong>Wait time: 15 minutes</strong></p>
  <p>🎁 <strong>Special: 20% off appetizers today!</strong></p>
  <p>Italian • $$ • 0.8 miles away</p>
</div>

<div style="border: 1px solid #ddd; padding: 15px; margin: 10px 0;">
  <h4>2. Pasta Palace ★ 4.3</h4>
  <p>📍 456 Elm St, New York, NY 10002</p>
  <p>🕒 <strong>Wait time: 25 minutes</strong></p>
  <p>🎁 No special offers today</p>
  <p>Italian • $$$ • 1.2 miles away</p>
</div>
```

```
<h3>Trending Near You</h3>
```

```
<div style="border: 1px solid #ddd; padding: 15px; margin: 10px 0;">
  <h4>6. Sushi House ★ 4.7</h4>
  <p>📍 789 Oak Ave, New York, NY 10003</p>
  <p>🕒 <strong>Wait time: 30 minutes</strong></p>
  <p>🎁 <strong>Special: Free dessert with any entree!</strong></p>
  <p>Japanese • $$$ • 0.5 miles away • 347 likes this week</p>
</div>
```

## Key Information Displayed

- ★ Yelp rating
- 📍 Full address
- 🕒 **Current wait-time** (NEW)
- 🎁 **Special offer** (NEW)
- Distance from user
- Price range
- Like count (for trending)

---

## 8. Key Assumptions and Design Decisions

### Wait-Time API Assumptions

- Each restaurant provides REST API endpoint
- Response time: < 1 second
- Rate limit: 1,000 requests/hour per API key
- Returns numeric wait-time in minutes
- Availability: 99% uptime
- Fallback: Show "N/A" if API unavailable
- Update frequency: Real-time, changes every 5-10 minutes

### Special Offers API Assumptions

- Daily offers released at midnight UTC

- Offers valid for 24 hours (expire next midnight)
- API returns text description (max 200 chars)
- Not all restaurants have offers every day
- Availability: 99.5% uptime
- No mid-day offer changes

## Caching Assumptions

- Wait-times acceptable if 5 minutes stale
- Special offers don't change during the day
- DynamoDB TTL cleanup: within 48 hours (AWS guarantee)
- Cache hit rate > 95% during steady state
- Cold start: First user may see "N/A" until cache warms

## Cost Assumptions

- External wait-time API: \$0.001 per call
- External offers API: \$0.0005 per call (less frequent)
- 5,000 restaurants in system
- 100,000 active users
- 10,000 recommendation requests/day

# 9. AWS Services Summary

All services from Assignment 1 (no new services added):

Service	Usage	Configuration for Question 2
S3	Frontend hosting	Same as Question 1
API Gateway	REST APIs	Same 4 endpoints
Lambda	Business logic	<b>8 functions</b> (was 6 in Q1)
Amazon Lex	Chatbot NLP	Same DiningSuggestionsIntent
SQS	Message queues	Same 2 queues (Q1, Q2)
DynamoDB	Data storage	<b>6 tables</b> (was 4 in Q1)

Service	Usage	Configuration for Question 2
ElasticSearch	Search engine	Same index, no changes
SES	Email delivery	Enhanced email template
CloudWatch	Monitoring & scheduling	<b>3 event rules</b> (was 1 in Q1)








**Total: Still only 9 AWS services** 

#### New Components:

- +2 Lambda functions (LF7, LF8)
- +2 DynamoDB tables (WaitTimeCache, SpecialOffers)
- +2 CloudWatch event rules

## 10. Conclusion

This solution successfully extends Question 1 with real-time wait-times and daily special offers integration, demonstrating:

-  **Same AWS Services** - No new infrastructure services needed
-  **Aggressive Caching** - 95%+ cache hit rate dramatically reduces costs
-  **Low Latency** - ~310ms total (only +50ms vs Question 1)
-  **Scalable** - Handles 100K users with concurrent external API calls
-  **Cost-Effective** - ~\$440/month (with intelligent API call optimization)
-  **Reliable** - Graceful degradation when external APIs fail
-  **Event-Driven** - Background jobs decouple user requests from slow APIs

#### Key Innovations:

1. **Cache-Aside Pattern:** User requests never wait for external APIs
2. **TTL-Based Expiration:** Automatic cache cleanup without manual intervention
3. **Multi-Layer Fallback:** System remains functional even with API outages
4. **Scheduled Background Jobs:** Expensive operations happen off peak hours

**The architecture proves that external API integration can be practical at scale through intelligent caching strategies!** 