

## AI Presence - Core Design [Updated Summary]

### Architecture

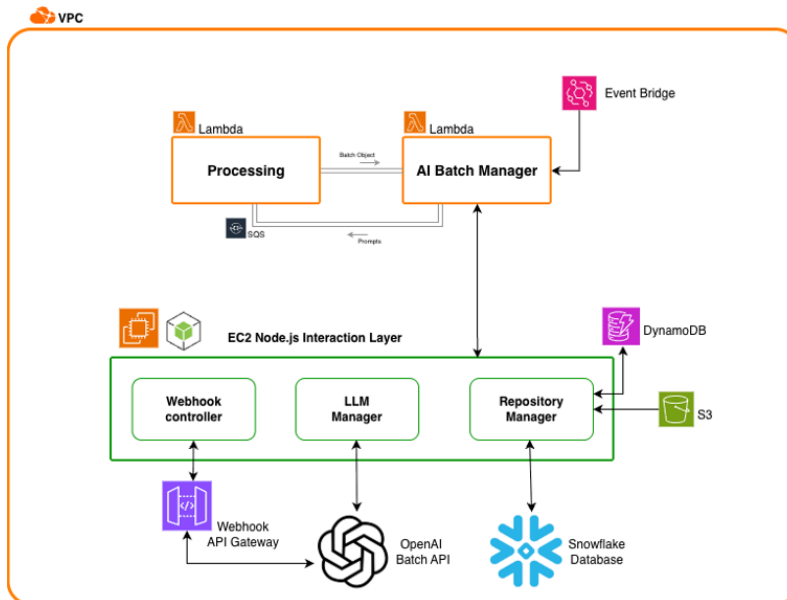
#### Description

A **serverless event-driven architecture**, processing Lambda which using a strategy patten, builds the relevant data pipeline on demand, coordinates closely with an EC2-hosted Node.js Communication Service, which provides orchestration, webhook control, LLM management, and repository handling. Communication between services remains asynchronous and decoupled.

#### Components

| Component                     | Purpose   |
|-------------------------------|---|
| VPC                           | Isolates all Lambda functions and EC2 service within secure network boundaries  |
| Processing Lambda             | Dynamic Processing Lambda that orchestrates all phases of the data flow using an internal Strategy Pattern and dynamic pipeline configuration. Constructs the api request at the single level for the batch manager       |
| AI Batch Manager              | Handle prompt-triggered batch processing (initialize, process per provider/phase, and finalize). Continuously update job entities across the workflow.  |
| EC2 Node.js Interaction Layer | Houses three core managers: Webhook controller (batch notification updates), LLM Manager interfaces with Model , Repository Manager manages data persistence to Snowflake, Presence Job on dynamoDB, config stored in S3. |
| DynamoDB                      | Provides persistent storage for storing data at prompt level to escort the lifecycle of a prompt (PresenceJob)  |
| EventBridge                   | Triggers a scheduled daily job.   |
| S3                            |   |
| API Gateway                   | Exposes webhook endpoints for external system callbacks (OpenAI API Gateway)  |

|                    |  |
|--------------------|--|
| Snowflake Database | External data warehouse for persistent storage of processed results.   |
| OpenAI Batch API   | External LLM service called by LLM Manager for AI/ML processing tasks. |



## Interaction layer

Docs:

[AI Presence - Interaction Layer - API](#)

## S3 configurations object

```

1  {
2    "open_ai": {
3      "presence_batch": {
4        "method": "string",
5        "system_prompt": "string",
6        "model": "string",
7        "max_res_tokens": number,
8        "url": "string",
9        "temperature": number
10     },
11     "analysis_batch": {
12       "method": "string",
13       "system_prompt": "string",
14       "model": "string",
15       "max_res_tokens": number,
16       "url": "string",
17       "temperature": number,
18       "output_schema": {json_schema} *
19     }
20   }
21 }

```

json schema\*

```

1 {
2   "format": {
3     "type": "json_schema",
4     "name": "brand_mentions_by_sentiment",
5     "schema": {
6       "type": "object",
7       "patternProperties": {
8         "^[\\w\\s-]+$": {
9           "type": "object",
10          "properties": {
11            "negative": { "type": "number" },
12            "neutral": { "type": "number" },
13            "positive": { "type": "number" }
14          },
15          "required": ["negative", "neutral", "positive"],
16          "additionalProperties": false
17        }
18      }
19    },
20    "strict": true
21  }
22 }
23

```

## Processing Lambda

The lambda utilizes a strategy and builder design patterns in order to costumly adjust the required data pipeline required for each stage of the process.

### prompt-object

```

1 {
2   "promptId" : "string",
3   "jobId"? : string,
4   "topic" : "string",
5   "promptText" : "string",
6   "batchId"? : "string",
7   "method" : "string",
8   "systemPrompt": "string",
9   "model": "string",
10  "llmProvider": "open_ai", // currently only one llm provider is
    supported
11  "phase": "initial" | "follow_up"
12  "maxResTokens": number,
13  "url": "string",
14  "temperature": number,
15  "vectorStoreId": "string",
16  "outputSchema"? : object {json_schema*}
17 }

```

### batch-object

```

1 {
2   "custom_id": "prompt_id",
3   "method": "method",
4   "url": "url",
5   "body": completions-request | responses-request
6 }

```

### completions-request

```

1 {
2   "model": "{ model }", // REQUIRED: model ID
3   "temperature": {temperature},
4   "max_tokens": number
5   "messages": [
6     {
7       "role": "system",
8       "content": "{system_prompt}."
9     },

```

```

10     {
11         "role": "user",
12         "content": "{{prompt_text}}!"
13     }
14 ]
15 }
16

```

## responses-request

```

1  {
2  "model": "{{ model }}",
3  "temperature": {{temperture}},
4  "max_output_tokens": number,
5  "text": {{ json_schema }},
6  "store": false,
7  "instructions": string {{brands-key}}
8  "input": [
9      {
10         "type" : "message", // always "message"
11         "role" : "system",
12         "content": {{system_prompt}}
13     },
14     {
15         "type" : "message", // always "message"
16         "role" : "user",
17         "content": {{prompt_text}}
18     }
19 ]
20 }

```

## Data

### The PresenceJob Entity

#### Overview

PresenceJob represents the **entire execution lifecycle** for a single prompt's daily AI Presence run.

#### DynamoDB Design

- **TTL:** 24 hours (automatic cleanup after each daily cycle).
- **Atomic updates:** Safe concurrent updates for webhooks and batch processing.
- **Update mechanism:**
  - Webhook events update job status atomically.
  - batches use separate fields to avoid overwriting.

```

1  {
2  id: string; // Unique ID (UUID)
3  topic: string; // Topic of analysis
4  promptId: string; // Associated prompt identifier
5  promptResponseId?: string // Prompt response Id
6  startTimestamp: string; // ISO start time
7  lastUpdateTimestamp?: string; // Filled at each job update
8  error?: string; // Any stage error
9  llmProvider: "open_ai" // currently a single provider
supported
phase: "initial" | "follow_up" // explicitly state job's phase
11
12  batches: [{
13      batchId: string; // Batch ID (OpenAI)
14      batchInputFileId: string;
15      batchOutputFileId?: string
16      batchStatus: 'QUEUED' | 'RUNNING' | 'COMPLETED' | 'FAILED';
17      tokenUsage?: {
18          inputTokens: number;
19          outputTokens: number;
20          totalTokens: number;
21      };
22  },
23  ...
24  ]
25

```

## Snowflake Database Design

### Overview

Our **Snowflake database** serves as the **source of truth** for historical visibility data, allowing us to query, aggregate, and visualize insights over time.

### PROMPT Table

| Column      | Type         | Description   |
|-------------|--------------|---|
| PROMPT_ID   | VARCHAR (PK) | Unique identifier for the prompt                    |
| TOPIC       | VARCHAR      | Topic prompt associated with                        |
| PROMPT_TEXT | TEXT         | The full prompt text                                |
| CREATED_AT  | TIMESTAMP    | When prompt was created                             |
| IS_ACTIVE   | BOOLEAN      | Enables prompt activation/deactivation for rotation |

### PROMPT\_RESPONSE Table

| Column        | Type         | Description                        |
|---------------|--------------|------------------------------------|
| PROMPT_ID     | VARCHAR      | Unique identifier for the prompt   |
| DATE          | Date         | Date                               |
| RESPONSE_ID   | VARCHAR (PK) | Unique identifier for the response |
| RESPONSE_TEXT | TEXT         | The full response.                 |

### PRESENCE\_SENTIMENT\_MENTIONS Table

| Column      | Type  | Description                        |
|-------------|---|------------------------------------|
| PROMPT_ID   | VARCHAR<br>(FK → PROMPT . PROMPT_ID )                 | The prompt used in this analysis   |
| RESPONSE_ID | VARCHAR (PK)<br>(FK → PROMPT_RESPONSE . RESPONSE_ID ) | Unique identifier for the response |
| MODEL       | VARCHAR   | model (e.g gpt-4.1-nano)           |

|                   |              |                                      |
|-------------------|--------------|--------------------------------------|
| BRANDKEY          | VARCHAR (PK) | Brand name analyzed                  |
| POSITIVE_MENTIONS | INTEGER      | Count of positive mentions           |
| NEUTRAL_MENTIONS  | INTEGER      | Count of neutral mentions            |
| NEGATIVE_MENTIONS | INTEGER      | Count of negative mentions for brand |
| DATE              | DATE         | Analysis date (daily granularity)    |
| TOTAL_TOKENS      | INTEGER      | Tokens consumed in LLM analysis      |

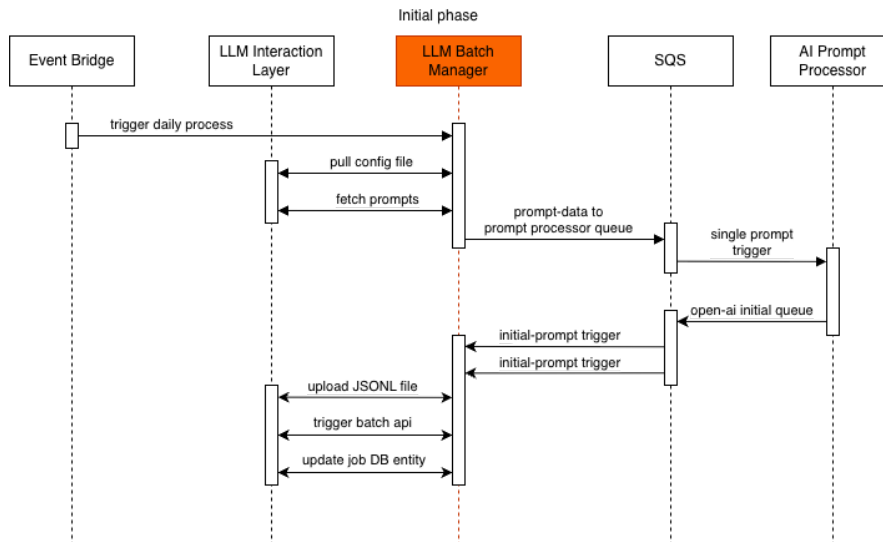
## 🔧 AI Batch Manager

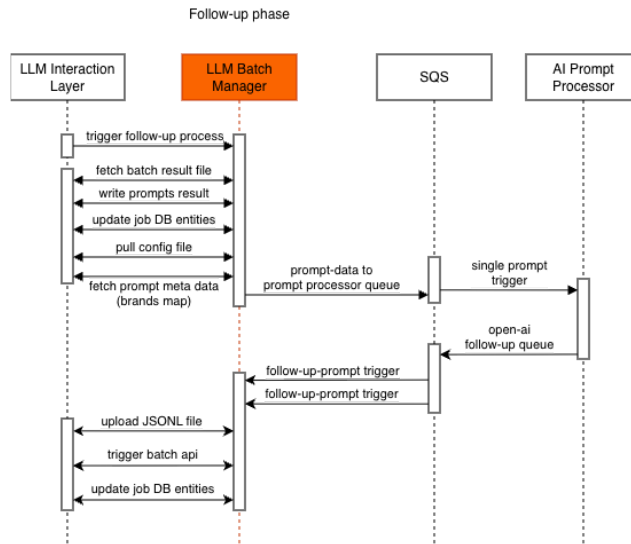
### Responsibilities Overview

- Prompt triggering** - Create config (per prompt) → trigger (using SQS) the AIPromptProcessor lambda
  - per initial/follow-up phase
- Batch initialization** - Consume SQS messages created by the AIPromptProcessor lambda → initialize batch process
  - per LLM provider → per initial/follow-up phase
- Batch Finalization**
- Throughout the flow update job entities as needed

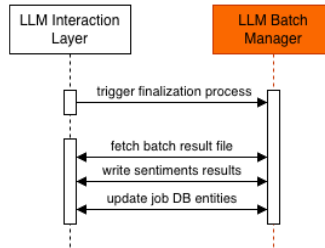
### Flows

\* Mind the following schemas neglects [ interaction layer ↔ ai prompt processor ] processes, for feature focusing & simplicity





Finalization phase



**Lambda ↔ SQS Limitations**

AWS Lambda has a built-in & non-configurable mechanism to Manage load balancing and functions invocations. This means that multiple “AI batch manager” functions can be triggered for a stream of SQS messages fired around at the same time.

**Implications**

ANY size of prompt count WILL be split and handled by multiple “AI Batch Manager” functions, leading to multiple batches.

**Short (?) term solution**

Limiting the “Maximum concurrency” of SQS triggers to “2” (2 being the minimum allowed by AWS).

Meaning:

10 prompts → 2 batches (each one with ~5)

1000 prompts → 2 batches (each one with ~500)

**AI Batch Manager ↔ SQS ( ↔ AI Prompt Processor ) - I/O specification**

Refer to [AI Presence - Core Design \[Updated Summary\] | Processing Lambda](#)

**EventBridge → AI Batch Manager - integration**

Daily trigger JSON input (from eventBridge)

```

1 {
2   eventSource: "event_bridge"
3   phase: "initial"
4 }
  
```

## AI Batch Manager ↔ Interaction Layer - I/O specification

**i** This section represents only the events for lambda invocation by ai-presence-interaction-layer. for full integration docs refer to [AI Presence - Interaction Layer - API | Update Job by ID](#)

Trigger AI batch manager (used by Interaction layer at Webhook invocations)

Usage notes: manual lambda invocation using *aws-sdk*

Invocation event:

```
1 {  
2   eventSource: "manual_invocation", // this exact string value  
3   phase: "follow_up" | "final", // these exact string values  
4   llmProvider: "open_ai" // this exact string value  
5   responseId: string // e.g. open-ai's output-file-id  
6   batchId: string // batchId  
7 }
```

Response (200): empty