

# SupportBot Multimodal Implementation Report

## Implementation Details, Evaluation Setup, and Examples

AI Agent (Cursor)

February 9, 2026

### Abstract

This report documents the implementation of multimodal (text + images) support in SupportBot. It describes how image attachments are persisted, how images are passed into Gemini model calls via Google's OpenAI-compatible endpoint, and what evaluation scripts are included in the repository. Quantitative results depend on the availability of decrypted Signal history/attachments and API credentials; this document focuses on reproducible behavior verifiable from the codebase.

## Contents

<b>1 Executive Summary</b>	<b>3</b>
1.1 Key Achievements . . . . .	3
1.2 Implementation Status . . . . .	3
<b>2 Algorithms (Current Implementation)</b>	<b>3</b>
2.1 Algorithm 1: Multimodal Message Ingestion . . . . .	4
2.2 Algorithm 2: Case Extraction with Validation . . . . .	6
2.3 Algorithm 3: Multimodal Response Pipeline . . . . .	9
<b>3 Examples (Illustrative)</b>	<b>11</b>
3.1 Example 1: Compatibility Question (synthetic) . . . . .	11
3.1.1 Raw Messages (Input) . . . . .	11
3.1.2 Extracted Case Block . . . . .	11
3.1.3 Structured Case (Output) . . . . .	11
3.1.4 Embedding & Storage . . . . .	11
3.2 Example 2: Multimodal Message with Screenshot (synthetic) . . . . .	12
3.2.1 Raw Messages (Input) . . . . .	12
3.2.2 Structured Case with Image Paths . . . . .	12
3.2.3 How Images Are Used . . . . .	12
<b>4 Solved Cases: Retrieval Introspection</b>	<b>13</b>
4.1 Example Query 1: Gimbal Control Issue . . . . .	13
4.1.1 User Question . . . . .	13
4.1.2 Stage 1: Semantic Search . . . . .	13
4.1.3 Stage 2: Image Loading . . . . .	13
4.1.4 Stage 3: LLM Decision . . . . .	14
4.2 Example Query 2: Configuration Parameter Question . . . . .	15
4.2.1 User Question . . . . .	15
4.2.2 Retrieval Results . . . . .	15
4.2.3 Bot Response . . . . .	15
4.3 Example Query 3: Multimodal Response (with Screenshot) . . . . .	16

4.3.1	User Question with Image . . . . .	16
4.3.2	Image Processing Pipeline . . . . .	16
4.3.3	Bot Response . . . . .	16
<b>5</b>	<b>Evaluation and Reproducibility</b>	<b>17</b>
5.1	Unit tests (offline) . . . . .	17
5.2	LLM-backed quality evaluation (requires <code>GOOGLE_API_KEY</code> ) . . . . .	17
5.3	Real-data evaluation (optional; requires decrypted Signal history) . . . . .	17
5.4	Notes on Signal Desktop encryption . . . . .	17
5.5	Embedding model note . . . . .	17
<b>6</b>	<b>Configuration and Limits</b>	<b>18</b>
6.1	Multimodal Settings . . . . .	18
6.2	Cost considerations . . . . .	18
<b>7</b>	<b>Conclusion</b>	<b>18</b>

# 1 Executive Summary

## 1.1 Key Achievements

Capability	Evidence in codebase
Persist image attachment paths	<code>raw_messages.image_paths_json</code> (see <code>signal-bot/app/db/schema.py</code> , <code>signal-bot/app/db/schema_mysql.py</code> )
Image-to-text extraction at ingestion	<code>signal-bot/app/ingestion.py</code> calling <code>LLMClient.image_to_text_json()</code>
Pass images into gate/response LLM calls	<code>signal-bot/app/jobs/worker.py</code> passes images into <code>LLMClient.decide_consider()</code> and <code>LLMClient.decide_and_respond()</code>
Store evidence image paths on cases	<code>signal-bot/app/jobs/worker.py:_collect_evidence_image_paths()</code> and DB field <code>cases.evidence_image_paths_json</code>
Surface KB evidence images at response time	<code>signal-bot/app/jobs/worker.py</code> loads images from retrieved case metadata <code>evidence_image_paths</code>

Table 1: Multimodal capabilities implemented (verifiable in code)

## 1.2 Implementation Status

All priority items from the proposed fix have been implemented:

- ✓ **P0:** Reject cases without `solution_summary` (High impact)
- ✓ **P1:** Pass images to `decide_and_respond()` (High impact)
- ✓ **P2:** Pass images to `decide_consider()` (Medium impact)
- ✓ **P3:** Store image paths in `raw_messages` (Enables P1/P2)
- ✓ **P4:** Include images in KB case evidence (Medium impact)

# 2 Algorithms (Current Implementation)

This section presents pseudoalgorithms that mirror the current multimodal implementation in the codebase.

## 2.1 Algorithm 1: Multimodal Message Ingestion

---

**Algorithm 1** Multimodal Message Ingestion — Preserves image paths for later use

---

```

1: procedure INGESTMESSAGE(msg_id, group_id, sender, ts, text, image_paths)
2:   content_text  $\leftarrow$  text
3:   context_text  $\leftarrow$  text
4:   stored_image_paths  $\leftarrow$  []
      ▷ NEW: Track valid image paths
5:
6:   for path in image_paths do
7:     img_path  $\leftarrow$  RESOLVEATTACHMENTPATH(path, settings.signal_bot_storage)
8:     img_path  $\leftarrow$  RESOLVE(img_path)
9:     if  $\neg$ img_path.EXISTS() then
10:      LOG.WARNING("Attachment missing: {path}")
11:      continue
12:    end if
13:
14:    stored_image_paths.APPEND(img_path) ▷ Store canonical path
15:
16:    ▷ Extract text/observations for searchability
17:    img_bytes  $\leftarrow$  READFILE(img_path)
18:    extraction  $\leftarrow$  LLM.IMAGETOTEXTJSON(img_bytes, context_text)
19:    content_text  $\leftarrow$  content_text + "[image]" + JSON(extraction)
20:    ▷ On extraction error, store placeholder JSON {observations: [], extracted_text: ""}
21:  end for
22:
23:  ▷ Store image paths alongside text
24:  INSERTRAWMESSAGE(msg_id, group_id, ts, sha256(sender)[: 16], content_text, stored_image_paths, reply_to)
25:
26:  ENQUEUEJOB(BUFFER_UPDATE, {group_id, msg_id})
27:  ENQUEUEJOB(MAYBE RESPOND, {group_id, msg_id})
28: end procedure

```

---

### Notable behavior:

- Image paths are stored in the database for later multimodal calls
- Relative attachment paths are resolved against **SIGNAL\_BOT\_STORAGE**
- Image-to-text extraction output is appended to message text for searchability



## 2.2 Algorithm 2: Case Extraction with Validation

---

**Algorithm 2** Case Extraction with Solution Validation — Reject solved cases without solutions

---

```

1: procedure HANDLEBUFFERUPDATE(group_id, msg_id)
2:   msg  $\leftarrow$  GETRAWMESSAGE(msg_id)
3:   line  $\leftarrow$  FORMATBUFFERLINE(msg)
4:   buffer  $\leftarrow$  GETBUFFER(group_id)
5:   buffer_new  $\leftarrow$  buffer + line
6:
7:   extract  $\leftarrow$  LLM.EXTRACTCASE(buffer_new)
8:   if extract.found then
9:     SETBUFFER(group_id, buffer_new)
10:    return
11:   end if
12:
13:   case  $\leftarrow$  LLM.MAKECASE(extract.case_block)
14:   if case.keep then
15:     SETBUFFER(group_id, extract.buffer_new)
16:     return
17:   end if
18:
19:                                $\triangleright$  Reject solved cases without solutions
20:   if case.status = "solved"  $\wedge$  case.solution_summary.STRIP() = "" then
21:     LOG.WARNING("Rejecting solved case without solution_summary")
22:     SETBUFFER(group_id, extract.buffer_new)
23:     return
24:   end if
25:
26:   case_id  $\leftarrow$  NEWUUID()
27:
28:                                $\triangleright$  Collect image paths from evidence messages
29:   evidence_image_paths  $\leftarrow$  COLLECTEVIDENCEIMAGES(case.evidence_ids)
30:
31:   INSERTCASE(case_id, group_id, case.*, evidence_image_paths)
32:
33:   doc_text  $\leftarrow$  JOINLINES(case.problem_title, case.problem_summary,
34:                           case.solution_summary, "tags: " + JOIN(case.tags))
35:   embedding  $\leftarrow$  LLM.EMBED(doc_text)
36:
37:                                $\triangleright$  Store image paths in metadata for retrieval
38:   CHROMA.UPSERT(case_id, doc_text, embedding,
39:                  {group_id, status, evidence_ids, evidence_image_pathsgroup_id, extract.buffer_new)
41: end procedure
42:
43: procedure COLLECTEVIDENCEIMAGES(evidence_ids)
44:   paths  $\leftarrow$  []
45:   for msg_id in evidence_ids do
46:     msg  $\leftarrow$  GETRAWMESSAGE(msg_id)
47:     if msg  $\neq$  null then
48:       for p in msg.image_paths do
49:         paths.APPEND(p)
50:       end for
51:     end if
52:   end for
53:   return paths

```

**Notable behavior:**

- Solved cases must have non-empty solution summaries
- Evidence image paths collected from raw messages
- Image paths stored in vector DB metadata for later retrieval



### 2.3 Algorithm 3: Multimodal Response Pipeline

---

**Algorithm 3** Multimodal Response Pipeline — Images at every decision point

---

```

1: procedure HANDLEMAYBERESPOND(group_id, msg_id)
2:   msg  $\leftarrow$  GETRAWMESSAGE(msg_id)                                 $\triangleright$  Now includes image_paths
3:   context  $\leftarrow$  GETLASTNMESSAGES(group_id, n)
4:
5:    $\triangleright$  Load images from current message for gate
6:   msg_images  $\leftarrow$  LOADIMAGES(msg.image_paths, max_gate, budget)
7:
8:   force  $\leftarrow$  MENTIONSBOT(msg.content_text)
9:   if  $\neg$ force then
10:     $\triangleright$  Gate sees images
11:    decision  $\leftarrow$  LLM.DECIDECONSIDER(msg.content_text, context, msg_images)
12:    if  $\neg$ decision.consider then
13:      return                                               $\triangleright$  Ignore greeting/noise
14:    end if
15:   end if
16:
17:   query_embedding  $\leftarrow$  LLM.EMBED(msg.content_text)
18:   retrieved  $\leftarrow$  CHROMA.RETRIEVE(group_id, query_embedding, k)
19:
20:    $\triangleright$  Collect images from retrieved KB cases
21:   kb_paths  $\leftarrow$  []
22:   for item in retrieved do
23:     paths  $\leftarrow$  item.metadata.evidence_image_paths
24:     kb_paths.EXTEND(paths[ $: \text{max\_per\_case}$ ])
25:   end for
26:   kb_paths  $\leftarrow$  kb_paths[ $: \text{max\_total\_kb}$ ]
27:
28:    $\triangleright$  Load KB images (respecting budget after msg images)
29:   remaining_budget  $\leftarrow$  MAX(budget - TOTALSIZE(msg_images), 0)
30:   kb_images  $\leftarrow$  LOADIMAGES(kb_paths, max_respond, remaining_budget)
31:
32:   all_images  $\leftarrow$  msg_images + kb_images
33:   all_images  $\leftarrow$  all_images[ $: \text{max\_images\_per\_respond}$ ]           $\triangleright$  Final cap
34:
35:   cases_json  $\leftarrow$  JSON(retrieved)
36:
37:    $\triangleright$  Responder sees all images
38:   resp  $\leftarrow$  LLM.DECIDEANDRESPOND(msg.content_text, context,
39:                                     cases_json, all_images)
40:
41:   if resp.respond then
42:     SIGNAL.SEND(group_id, resp.text)
43:   end if
44: end procedure
45: procedure LOADIMAGES(paths, max_count, budget_bytes)
46:   images  $\leftarrow$  []
47:   total  $\leftarrow$  0
48:   for p in paths do
49:     if  $|\text{images}| \geq \text{max\_count}$  then           9
50:       break
51:     end if
52:     data  $\leftarrow$  READFILE(p)
53:     image  $\leftarrow$  [ data ]
54:
```

**Notable behavior:**

- **P2:** Gate stage receives images from user message
- **P1:** Responder receives images from both user message and KB evidence
- **P4:** Evidence images retrieved from case metadata
- Image budgets cap multimodal payload size (`MAX_IMAGE_SIZE_BYTES` and `MAX_TOTAL_IMAGE_BYTES`)

## 3 Examples (Illustrative)

This section provides illustrative, synthetic examples to show how messages are transformed into structured cases and how images flow through the pipeline. Message IDs, timestamps, similarity values, and file paths shown below are placeholders and will differ in real runs.

### 3.1 Example 1: Compatibility Question (synthetic)

#### 3.1.1 Raw Messages (Input)

User:

```
Is device MODEL_X supported? Which target should I select when flashing?
```

Support:

```
Yes. Use target TARGET_Y. The tooling will detect it as TARGET_Y.
```

#### 3.1.2 Extracted Case Block

CASE BLOCK:

```
problem: Device compatibility and flashing target selection
```

```
evidence: [msg_id_1, msg_id_2]
```

```
status: solved
```

#### 3.1.3 Structured Case (Output)

Field	Value
case_id	c4f2a891-...
status	solved
problem_title	Support for device MODEL_X
problem_summary	User asks whether MODEL_X is supported and what target to select when flashing.
solution_summary	Support confirms compatibility and recommends selecting target TARGET_Y in the flashing tool.
tags	compatibility, firmware, flashing, target selection
evidence_ids	[msg_id_1, msg_id_2]
evidence_image_paths	(no images in this case)

Table 2: Structured case with metadata

#### 3.1.4 Embedding & Storage

- **Document text:** Built from title + summaries + tags (see `signal-bot/app/jobs/worker.py`)
- **Embedding:** Vector produced by the configured `EMBEDDING_MODEL`
- **Vector DB:** Stored in ChromaDB with metadata: {`group_id`, `status`, `evidence_ids`, `evidence_image_paths`}

## 3.2 Example 2: Multimodal Message with Screenshot (synthetic)

### 3.2.1 Raw Messages (Input)

User:

Please help. I see an error when trying to complete an action. What does it mean?  
[image: <attachment.png>]

Support:

Which mode/step are you in when the error appears?

User:

It happens in MODE\\_A. In MODE\\_B it works normally.

Support:

Please share logs/config. The screenshot may include an error code we can diagnose.

### 3.2.2 Structured Case with Image Paths

Field	Value
problem_title	Error shown in screenshot during MODE_A
problem_summary	User cannot complete an action in MODE_A; screenshot contains an error message/code. MODE_B works.
solution_summary	Gather logs/config and troubleshoot based on the exact error text visible in the screenshot and retrieved evidence cases.
evidence_image_paths	<code>["&lt;abs_path_to_attachment.png&gt;"]</code>

### 3.2.3 How Images Are Used

At ingestion:

- Image extracted to text: {observations: ["..."], extracted\_text: "...”}
- Image path stored: <abs\_path\_to\_attachment.png>

At retrieval (when user asks similar question):

1. User query: "What does this screenshot error mean?"
2. System retrieves relevant case(s) via semantic similarity search
3. Loads evidence image(s) from disk (bounded by size/budget limits)
4. Passes image(s) + retrieved case text to the LLM for response generation

## 4 Solved Cases: Retrieval Introspection

This section demonstrates how the bot retrieves and reasons about cases when answering user questions, with full introspection into the retrieval pipeline.

### 4.1 Example Query 1: Gimbal Control Issue

#### 4.1.1 User Question

User: Need to control FEATURE\\_X but it doesn't work on this configuration.  
It looks like DEVICE\\_Y lacks MODE\\_Z. RC passthrough alone is not enough.

#### 4.1.2 Stage 1: Semantic Search

**Query embedding:** Generated from user question

**Search parameters:**

- `group_id: <group_id>`
- `k: RETRIEVE_TOP_K` (default: 5)
- `embedding_model: EMBEDDING_MODEL` (configurable)

**Retrieved cases (ranked by similarity score):**

Rank	Case Title
1	FEATURE_X control troubleshooting
2	Output/servo mode configuration
3	Missing MODE_Z support on DEVICE_Y
4	RC passthrough configuration
5	Firmware option enabling FEATURE_X

Table 3: Top-5 retrieved cases

#### 4.1.3 Stage 2: Image Loading

**For each retrieved case:**

- Case 1: `evidence_image_paths = []` (no images)
- Case 2: `evidence_image_paths = []` (no images)
- Case 3: `evidence_image_paths = []` (no images)
- Case 4: `evidence_image_paths = []` (no images)
- Case 5: `evidence_image_paths = []` (no images)

**Total images loaded: <n>**

**Total budget used: <bytes> / MAX\_TOTAL\_IMAGE\_BYTES**

#### 4.1.4 Stage 3: LLM Decision

**Input to LLM:**

- User message (text)
- Context: last 40 messages from group
- Retrieved cases: JSON with top-5 cases
- Images: 0 message images + 0 KB images = 0 total

**LLM output:**

```
{  
  "respond": true,  
  "text": "<answer based on retrieved case evidence>",  
  "citations": ["case:<case_id>"]  
}
```

## 4.2 Example Query 2: Configuration Parameter Question

### 4.2.1 User Question

User: Good day. Analyzing parameter "PARAM\\_X": Value: V. If the system error exceeds this threshold, a failsafe will trigger. Is this too low? What are the recommendations?

### 4.2.2 Retrieval Results

Rank	Case Title	Status
1	PARAM_X configuration	solved
2	Failsafe parameters	solved
3	Sensor/estimation troubleshooting	solved

**Top case evidence:**

- **Problem:** User asks about a configuration parameter and whether the current value is appropriate
- **Solution:** <solution summary from retrieved case evidence>
- **Evidence IDs:** [msg\_id\_42, msg\_id\_43, msg\_id\_44]
- **Images:** No images in evidence

### 4.2.3 Bot Response

<answer based on retrieved case evidence>

Ref: case:<case\_id>...

## 4.3 Example Query 3: Multimodal Response (with Screenshot)

### 4.3.1 User Question with Image

User: Look at what the ground control app shows. Don't understand what kind of error this is. [image: <attachment.png>]

### 4.3.2 Image Processing Pipeline

#### Step 1: Message ingestion

- Image path stored in DB: <abs\_path\_to\_attachment.png>
- Image-to-text JSON appended to content\_text: {"observations": ["..."], "extracted\_text": "..."}

#### Step 2: Gate decision (decide\_consider)

- Message text: "Look at what the ground control app shows..."
- **Images loaded:** up to MAX\_IMAGES\_PER\_GATE (subject to size/budget limits)
- LLM sees: text + the attached image(s)
- Decision: consider=<true/false>

#### Step 3: Semantic retrieval

- Query embedding from: <user message text>
- Top retrieved case: <case title>
- Case may include evidence\_image\_paths: ["<abs\_path\_to\_evidence.png>"]

#### Step 4: Response generation (decide\_and\_respond)

- Message images: <n>
- KB images: <n>
- **Total images:** capped by MAX\_IMAGES\_PER\_RESPOND
- **Total bytes:** capped by MAX\_TOTAL\_IMAGE\_BYTES
- LLM sees: user image(s) + retrieved case text (+ optional KB evidence image(s))

### 4.3.3 Bot Response

<answer based on retrieved case evidence>

Ref: case:<case\_id>...

#### Why multimodal can help (conceptually):

- Without images: the model only sees text, which may be ambiguous
- With images: the model can directly see UI text/error codes and other visual context
- With KB evidence images: the model can compare the user's screenshot to historical evidence

## 5 Evaluation and Reproducibility

This repository contains tests and evaluation scripts. Most evaluation outputs are written under `test/data/` and are intentionally gitignored; quantitative metrics therefore depend on the local environment, available decrypted Signal history/attachments, and API credentials.

### 5.1 Unit tests (offline)

- Run: `pytest -v`
- Core coverage includes ingestion, buffer/case extraction, Chroma integration, and response gating.

### 5.2 LLM-backed quality evaluation (requires `GOOGLE_API_KEY`)

- Tests: `test/test_quality_eval.py`
- The judge uses Gemini via Google's OpenAI-compatible endpoint.

### 5.3 Real-data evaluation (optional; requires decrypted Signal history)

- Mine cases: `python test/mine_real_cases.py`
- Run eval: `python test/run_real_quality_eval.py`
- Streaming eval: `python test/run_streaming_eval.py`
- Image-to-text demo: `python test/run_image_to_text_demo.py`

### 5.4 Notes on Signal Desktop encryption

Signal Desktop backups may require Windows DPAPI decryption under the same Windows user account that created the backup (see `test/results.md` for a documented example).

### 5.5 Embedding model note

- Application default: `EMBEDDING_MODEL=text-embedding-004` (see `signal-bot/app/config.py`).
- The real-eval script may override to `gemini-embedding-001` for compatibility with the OpenAI-style Gemini embeddings endpoint (see `test/run_real_quality_eval.py`).

## 6 Configuration and Limits

### 6.1 Multimodal Settings

Parameter	Value	Purpose
MAX_IMAGES_PER_GATE	3	Limit images sent to gate decision
MAX_IMAGES_PER_RESPOND	5	Limit total images in response call
MAX_KB_IMAGES_PER_CASE	2	Limit evidence images per retrieved case
MAX_IMAGE_SIZE_BYTES	5,000,000	Skip images > 5,000,000 bytes
MAX_TOTAL_IMAGE_BYTES	20,000,000	Total budget per response (20,000,000 bytes)

Table 4: Image budget limits (caps multimodal payload size)

### 6.2 Cost considerations

The main cost drivers are LLM calls during ingestion and response:

- **Ingestion:** optional image-to-text extraction per attachment (`image_to_text_json`)
- **Retrieval:** embeddings for case documents and user queries (`embed`)
- **Response:** gate (`decide_consider`) and responder (`decide_and_respond`) chat calls, optionally with images

Actual costs depend on model selection and provider pricing. The implementation enforces strict caps on image count and total bytes to bound multimodal payload size (Table 4).

## 7 Conclusion

This implementation adds end-to-end multimodal plumbing:

1. **Reject low-quality cases (P0):** Reject `status=solved` cases without `solution_summary`
2. **Preserve image references (P3):** Store attachment paths in `raw_messages.image_paths_json`
3. **Use images in decisions and responses (P1, P2):** Pass images into `decide_consider` and `decide_and_respond`
4. **Carry evidence images through retrieval (P4):** Store and retrieve `evidence_image_paths` via Chroma metadata

**Measuring impact:** Use the scripts in Section 5 to run evaluations in an environment with decrypted data and valid API credentials.

**Next steps:**

- Deploy to production and monitor real-world performance
- Gather user feedback on response quality
- Fine-tune retrieval thresholds based on precision/recall metrics
- Consider adding image captioning for better searchability