

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

1 Executive Summary	3
1.1 Key Achievements	3
1.2 Implementation Status	3
2 Algorithms (Current Implementation)	3
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
3 Examples (Concrete Cases from Unit Tests & Evaluation Dataset)	11
3.1 Example 1: Login Problem (from unit test fixture)	11
3.1.1 Raw Messages (Input)	11
3.1.2 Structured Case (Output)	11
3.1.3 Embedding & Storage	11
3.2 Example 2: Flight Controller Error with Screenshot (from real evaluation dataset)	13
3.2.1 Raw Messages (Input)	13
3.2.2 Structured Case with Image Paths	13
3.2.3 How Images Are Used	13
4 Solved Cases: Retrieval Introspection	15
4.1 Example Query 1: Video Playback Issue	15
4.1.1 User Question	15
4.1.2 Stage 1: Semantic Search	15
4.1.3 Stage 2: Image Loading	15
4.1.4 Stage 3: LLM Decision	16
4.2 Example Query 2: Flight Controller Error with Multimodal Evidence	17
4.2.1 User Question with Image	17
4.2.2 Retrieval Results	17
4.2.3 Multimodal Response Pipeline	17

5 Evaluation and Reproducibility	18
5.1 Real Evaluation Dataset (400/100 Ukrainian Messages)	18
5.1.1 Dataset Composition	18
5.1.2 Evaluation Message Labels	18
5.2 Evaluation Results (Baseline)	18
5.3 Reproduction Instructions	19
5.4 Notes on Signal Desktop encryption	19
5.5 Embedding model note	19
6 Configuration and Limits	20
6.1 Multimodal Settings	20
6.2 Cost considerations	20
7 Conclusion	20

1 Executive Summary

1.1 Key Achievements

Capability	Evidence in codebase
Persist image attachment paths	<code>raw_messages.image_paths_json</code> <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> <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_pathsend procedure
41:
42: procedure COLLECTEVIDENCEIMAGES(evidence_ids)
43:   paths  $\leftarrow$  []
44:   for msg_id in evidence_ids do
45:     msg  $\leftarrow$  GETRAWMESSAGE(msg_id)
46:     if msg  $\neq$  null then
47:       for p in msg.image_paths do
48:         paths.APPEND(p)
49:       end for
50:     end if
51:   end for
52:   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[: max_per_case])
25:   end for
26:   kb_paths  $\leftarrow$  kb_paths[: 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[: 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  $|images| \geq max\_count$  then           9
50:       break
51:     end if
52:     data  $\leftarrow$  READFILE(p)
53:     images. $\leftarrow$  [data]

```

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 (Concrete Cases from Unit Tests & Evaluation Dataset)

This section provides concrete examples from the Ukrainian tech support fixture (`test/conftest.py:STABX_SU`) and real evaluation data (`test/data/streaming_eval/`). These examples demonstrate how real-world messages are transformed into structured cases.

3.1 Example 1: Login Problem (from unit test fixture)

3.1.1 Raw Messages (Input)

Source: `test/conftest.py`, Case 1 (Lines 387–393)

Group:

Timestamps: 1707400000000 – 1707400360000

```
user1 (ts=1707400000000):
! , ,
```

```
support1 (ts=1707400060000):
! cookies.
Caps Lock
```

```
user1 (ts=1707400120000):
,
```

```
support1 (ts=1707400180000):
.
```

```
user1 (ts=1707400300000):
, ! !
```

```
support1 (ts=1707400360000):
! -
```

3.1.2 Structured Case (Output)

Field	Value
case_id	<generated-uuid>
status	solved
problem_title	
problem_summary	, , . .
solution_summary	. .
tags	login, password, cache, recovery, personal-cabinet
evidence_ids	[msg_1707400000000, msg_1707400060000, ..., msg_1707400360000]
evidence_image_paths	[] (no images in this case)

Table 2: Structured case: login problem solved via password reset

3.1.3 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: Flight Controller Error with Screenshot (from real evaluation dataset)

3.2.1 Raw Messages (Input)

Source: test/data/streaming_eval/eval_messages_labeled.json, messages 1–4

Group: 019b5084-b6b0-7009-89a5-7e41f3418f98 ()

Timestamps: 1770285647836 – 1770293731770

Label: answer (requires technical response)

User (6928c2c3-1440-4215-98cf-6d6981c0d9c7) :

, ?

"PreArm: Internal Error 0x8000"

[ATTACHMENT image/png size=26467]

Support (85c10856-218e-4a35-bb63-53febaf61bf3) :

,

User (798dea6a-7d5e-44e1-be65-8e0a88b273b3) :

USB ?

Support (230003f4-75fc-4f20-ba0e-97aef2cc3c95) :

, .
,
. .

[ATTACHMENT image/jpeg size=25752]

3.2.2 Structured Case with Image Paths

Field	Value
problem_title	PreArm: Internal Error 0x8000 -
problem_summary	"PreArm: Internal Error 0x8000". . .
solution_summary	: (1) Matek , (2) USB, (3) - . . - . .
tags	flight-controller, prearm-error, internal-error, reboot, hardware-failure, matek
evidence_image_paths	["26/26c446716711fe8172591e0a539bfd9a97b2..."] , ["57/57c87921818f13999f4ab0fba6611ca70a11..."]

Table 3: Multimodal case with image evidence from real evaluation dataset

3.2.3 How Images Are Used

At ingestion:

- First image (26467 bytes PNG): extracted to text via LLM.ImageToTextJSON()
- Second image (25752 bytes JPEG): extracted to text
- Paths stored: 26/26c446716711fe8172..., 57/57c87921818f13...
- Content text includes: [image]{observations: [...], extracted_text: "PreArm: Internal Error 0x8000"}

At retrieval (when similar question asked):

1. User query: " Internal Error "
2. System retrieves this case via semantic similarity
3. Loads image(s) from `evidence_image_paths` (bounded by `MAX_KB_IMAGES_PER_CASE=2`)
4. Passes images + case text to `LLM.DecideAndRespond()` for visual context
5. Bot can reference the specific error code visible in screenshot

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 using real Ukrainian examples.

4.1 Example Query 1: Video Playback Issue

4.1.1 User Question

Source: Based on `test/conftest.py`, Case 2 (Video not playing)

User: , ,

4.1.2 Stage 1: Semantic Search

Query embedding: Generated from user question

Search parameters:

- `group_id: stabx-academy-support-group-123`
- `k: 5` (RETRIEVE_TOP_K default)
- `embedding_model: text-embedding-004`

Retrieved cases (ranked by similarity score):

Rank	Case Title (Ukrainian)
1	
2	
3	
4	
5	

Table 4: Top-5 retrieved cases from Ukrainian support chat

4.1.3 Stage 2: Image Loading

For each retrieved case:

- Case 1 (Video loading): `evidence_image_paths = []` (no images)
- Case 2 (Login issue): `evidence_image_paths = []` (no images)
- Case 3 (Payment): `evidence_image_paths = []` (no images)
- Case 4 (Mobile app): `evidence_image_paths = []` (no images)
- Case 5 (Progress lost): `evidence_image_paths = []` (no images)

Total images loaded: 0

Total budget used: 0 bytes / MAX_TOTAL_IMAGE_BYTES

4.1.4 Stage 3: LLM Decision

Input to LLM:

- User message: " , , "
- Context: last 40 messages from group
- Retrieved cases: JSON with case 1 solution: " Chrome Edge. Firefox "
- Images: 0 message images + 0 KB images = 0 total

LLM output:

```
{  
  "respond": true,  
  "text": " ! ? Firefox  
  
Chrome Edge - .",  
  "citations": ["case:<uuid-video-playback>"]  
}
```

4.2 Example Query 2: Flight Controller Error with Multimodal Evidence

4.2.1 User Question with Image

Source: Real evaluation dataset (message idx 8)

User: , ,
, [ATTACHMENT image/png size=169525]

4.2.2 Retrieval Results

Rank	Case Title	Status
1	PreArm: Internal Error 0x8000 -	solved
2	GPS	solved
3	Matek	solved

Table 5: Top-3 retrieved cases (from real 400/100 evaluation dataset)

Top case evidence:

- **Problem:** Internal Error 0x8000
- **Solution:** (1) Matek, (2) USB, (3)
- **Evidence IDs:** [msg_1770285647836, msg_1770286098545, msg_1770293731770]
- **Evidence Images:** ["26/26c44671...", "57/57c87921..."] (2 images, 52KB total)

4.2.3 Multimodal Response Pipeline

Images loaded:

- User message image: b3/b30d1e93867d3... (169525 bytes PNG)
- KB evidence images: 2 images from top case (52KB total)
- Total: 3 images, 221KB < MAX_TOTAL_IMAGE_BYTES (20MB)

Bot response (with citation):

,
:
1. SD- (Matek H743)
2. USB
3. -

,
USB . -
.

Ref: case:<uuid-internal-error-0x8000>

Trust features (implemented):

- ✓ Bot references at least one concrete solution from KB (case uuid cited)
- ✓ Bot quotes/mentions the original asker (via Signal quote feature in signal-cli)

5 Evaluation and Reproducibility

5.1 Real Evaluation Dataset (400/100 Ukrainian Messages)

Dataset source: test/data/streaming_eval/
Group: (019b5084-b6b0-7009-89a5-7e41f3418f98)
Created: 2026-02-09 10:58:28
Last evaluated: 2026-02-09 11:27:25

5.1.1 Dataset Composition

Component	Count	Purpose
Total messages used	500	Source material from real Signal group
Context messages	400	Build knowledge base and provide chat history
Evaluation messages	75	Test bot's question answering capability
Knowledge base cases	28	Extracted and embedded solved cases

Table 6: Real evaluation dataset structure (from `dataset_meta.json`)

5.1.2 Evaluation Message Labels

The 75 evaluation messages were labeled by LLM (`gemini-2.5-flash-lite`) into three categories:

Label	Count	Meaning
<code>answer</code>	23	Technical question requiring bot response
<code>contains_answer</code>	21	Message contains solution to previous question
<code>ignore</code>	31	Chatter, greetings, or acknowledgments (bot should not respond)

Table 7: Evaluation message label distribution

5.2 Evaluation Results (Baseline)

Source: test/data/streaming_eval/eval_summary.json

Label	N	Pass Rate	Avg Score	Respond Rate
<code>answer</code>	23	8.7%	0.96/10	13%
<code>contains_answer</code>	21	81.0%	8.10/10	19%
<code>ignore</code>	31	96.8%	9.68/10	3.2%
Overall	75	65.3%	6.56/10	—

Table 8: Baseline evaluation results (before multimodal + trust improvements)

Key insights:

- Bot correctly ignores most chatter (96.8% pass rate on `ignore` messages)
- Bot struggles with technical questions (`answer` label: only 8.7% pass rate, 0.96/10 avg score)

- Low response rate (13%) on questions requiring answers indicates overly conservative gating
- Overall pass rate 65.3% is dominated by correct "ignore" behavior rather than helpful answers

5.3 Reproduction Instructions

Prerequisites:

- Python 3.11+
- GOOGLE_API_KEY environment variable
- Decrypted Signal history/attachments (optional for real-data eval)

Unit tests (offline):

- Run: `pytest -v`
- Core coverage includes ingestion, buffer/case extraction, Chroma integration, and response gating.
- Uses synthetic Ukrainian fixtures from `test/conftest.py:STABX_SUPPORT_CHAT`

LLM-backed quality evaluation:

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

Real-data evaluation (requires decrypted Signal history):

- Prepare dataset: `python test/prepare_streaming_eval_dataset.py`
- Run evaluation: `python test/run_streaming_eval.py`
- Mine cases: `python test/mine_real_cases.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 9: 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 9).

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