

Fix Summary: Enhanced Chat 500 Error Resolution

Date: October 7, 2025

Issue: 500 Internal Server Error on `/api/ai/enhanced-chat` endpoint

Status:  RESOLVED

PR: [#123 - Enhanced Chat Streaming Support and 500 Error Resolution](https://github.com/dfulton-thebar/Sports-Bar-TV-Controller/pull/123) (<https://github.com/dfulton-thebar/Sports-Bar-TV-Controller/pull/123>)

Executive Summary

The chatbot was failing with a 500 Internal Server Error due to a missing `ai-knowledge-base.json` file. This caused silent failures during context building and extremely long response times (60+ seconds). The fix involved:

1. **Building the missing knowledge base** using the existing build script
2. **Adding streaming support** for real-time feedback and better UX
3. **Maintaining backward compatibility** with non-streaming mode

The chatbot now works correctly with both streaming and non-streaming modes, providing immediate feedback to users and successfully loading context from documentation and codebase.

Root Cause Analysis

Primary Issue: Missing Knowledge Base File

Location: `data/ai-knowledge-base.json`

Impact:

- The `loadKnowledgeBase()` function in `src/lib/ai-knowledge.ts` threw an error when the file was missing
- This error was caught in `buildEnhancedContext()` but only logged, not properly handled
- The endpoint continued processing without context, resulting in:
 - Very long response times (60+ seconds)
 - Poor quality responses without documentation context
 - Silent failures that were hard to debug

Why It Happened:

- The knowledge base file is in `.gitignore` (correctly, as it's a generated file)
- The file must be built using `npm run build-knowledge-base` after deployment
- This step was not performed after the recent optimizations were deployed

Secondary Issue: No Streaming Support

Impact:

- Users had to wait 60+ seconds with no feedback

- No indication that the system was working
- Poor user experience during AI processing

Solution Implemented

1. Built the Knowledge Base

Command: `npm run build-knowledge-base`

Results:

✓ Knowledge Base Built Successfully!



Statistics:

- Total Document Chunks: 2,754
- Total Files: 584
- PDF Chunks: 373
- Markdown Chunks: 555
- Code Chunks: 1,826
- Total Characters: 4,845,180
- File Size: 6.1MB
- Saved to: data/ai-knowledge-base.json

What It Does:

- Indexes all documentation files (PDFs, Markdown) from the `docs/` directory
- Indexes relevant codebase files (TypeScript, JavaScript, React components)
- Creates searchable chunks with metadata
- Enables context-aware AI responses

2. Added Streaming Support

File Modified: `src/app/api/ai/enhanced-chat/route.ts`

Changes:

- Added `handleStreamingResponse()` function for Server-Sent Events (SSE)
- Refactored existing code into `handleNonStreamingResponse()`
- Added status updates during context building
- Streams tokens in real-time as they're generated by Ollama
- Default behavior is now streaming (`stream: true`)

Benefits:

- Immediate feedback to users ("Building context...", "Generating response...")
- Real-time token streaming for better perceived performance
- Users can see the response being generated
- Significantly improved user experience

3. Maintained Backward Compatibility

Non-Streaming Mode: Available with `stream: false` parameter

- Returns complete response in JSON format
 - Compatible with existing clients
 - Useful for programmatic access
-

Testing Results

Test 1: Streaming Mode (Default)

Request:

```
curl -N -X POST http://localhost:3000/api/ai/enhanced-chat \
  -H "Content-Type: application/json" \
  -d '{"message": "What is CEC control?", "stream": true}'
```

Response (Server-Sent Events):

```
data: {"type": "status", "message": "Building context..."}
data: {"type": "context", "usedContext": true, "contextError": null}
data: {"type": "status", "message": "Generating response..."}
data: {"type": "token", "content": "CEC"}
data: {"type": "token", "content": " stands"}
data: {"type": "token", "content": " for"}
data: {"type": "token", "content": " Consumer"}
data: {"type": "token", "content": " Electronics"}
data: {"type": "token", "content": " Control"}
...
data: {"type": "done", "model": "llama3.2:3b", "usedContext": true, ...}
```

Results:

- ☒ Immediate status updates
- ☒ Real-time token streaming
- ☒ Context successfully loaded (usedContext: true)
- ☒ Response time: ~30 seconds with streaming feedback
- ☒ Content-Type: text/event-stream

Test 2: Non-Streaming Mode (Backward Compatible)

Request:





```
curl -X POST http://localhost:3000/api/ai/enhanced-chat \
  -H "Content-Type: application/json" \
  -d '{"message": "What is the AI assistant?", "stream": false}'
```

Response (JSON):

```
{
  "response":
    "The AI assistant in this context appears to be a machine learning-based system that
    analyzes floor plans and provides suggestions for TV placements...",
  "model": "llama3.2:3b",
  "usedContext": true,
  "usedCodebase": true,
  "usedKnowledge": true,
  "contextError": null
}
```

Results:

- ☒ Complete response returned





-  Context successfully loaded (`usedContext: true`)
-  Backward compatible with existing clients
-  Response time: ~55 seconds (expected for non-streaming)
-  Content-Type: `application/json`

Test 3: Without Context (Fast Mode)

Request:


```
curl -N -X POST http://localhost:3000/api/ai/enhanced-chat \
  -H "Content-Type: application/json" \
  -d '{"message": "Hello", "stream": true, "useKnowledge": false, "useCodebase": false}'
```

Results:





-  Immediate response (no context building delay)
-  Streaming works correctly
-  Response time: ~2 seconds
-  Useful for simple queries that don't need context

Performance Comparison

Before Fix

Metric	Value
Status Code	500 Internal Server Error
Response Time	60+ seconds (when it worked)
User Feedback	None (black box)
Context Loading	Failed silently
User Experience	 Very Poor

After Fix

Metric	Streaming Mode	Non-Streaming Mode
Status Code	200 OK	200 OK
Response Time	~30s with feedback	~55s
User Feedback	Real-time status & tokens	Complete response
Context Loading	 Success	 Success
User Experience	 Excellent	 Good

Deployment Instructions

For Production Deployment

1. **Pull the latest changes:**

```
bash
cd ~/Sports-Bar-TV-Controller
git pull origin fix/enhanced-chat-streaming-and-knowledge-base
```

2. **Build the knowledge base (REQUIRED):**

```
bash
npm run build-knowledge-base
```

⚠ This step is critical - the chatbot will not work without it!

3. **Rebuild the application:**

```
bash
npm run build
```

4. **Restart the server:**

```
bash
npm start
# or if using PM2:
pm2 restart sportsbar-assistant
```

For Development

1. **Pull the latest changes:**

```
bash
git pull origin fix/enhanced-chat-streaming-and-knowledge-base
```

2. **Build the knowledge base:**

```
bash
npm run build-knowledge-base
```

3. **Start development server:**

```
bash
npm run dev
```

Important Notes

Knowledge Base Maintenance

⚠ **Critical:** The `data/ai-knowledge-base.json` file is in `.gitignore` and must be built on each deployment.

When to Rebuild:

- After every deployment
- When documentation files are updated
- When significant codebase changes are made
- If the chatbot starts giving outdated information

How to Rebuild:

```
npm run build-knowledge-base
```

Streaming vs Non-Streaming

Use Streaming (Default):

- For interactive chatbot UI
- When users need immediate feedback
- For better perceived performance

Use Non-Streaming:

- For programmatic API access
- When you need the complete response at once
- For backward compatibility with existing clients

How to Control:

```
// Streaming (default)
fetch('/api/ai/enhanced-chat', {
  method: 'POST',
  body: JSON.stringify({ message: 'Hello', stream: true })
})

// Non-streaming
fetch('/api/ai/enhanced-chat', {
  method: 'POST',
  body: JSON.stringify({ message: 'Hello', stream: false })
})
```

Code Changes

Modified Files

1. `src/app/api/ai/enhanced-chat/route.ts`
 - Added streaming support with SSE
 - Refactored into `handleStreamingResponse()` and `handleNonStreamingResponse()`
 - Added status updates during processing
 - Improved error handling

Generated Files (Not in Git)

1. `data/ai-knowledge-base.json` (6.1MB)
 - Generated by `npm run build-knowledge-base`
 - Contains indexed documentation and codebase
 - Must be rebuilt after deployment
-

Verification Steps

1. Check Knowledge Base Exists

```
ls -lh data/ai-knowledge-base.json
# Should show: -rw-r--r-- 1 user user 6.1M Oct 7 20:59 data/ai-knowledge-base.json
```

2. Test Streaming Endpoint

```
curl -N -X POST http://localhost:3000/api/ai/enhanced-chat \
-H "Content-Type: application/json" \
-d '{"message": "test", "stream": true}'
```

Expected: Should see `data:` events streaming in real-time

3. Test Non-Streaming Endpoint

```
curl -X POST http://localhost:3000/api/ai/enhanced-chat \
-H "Content-Type: application/json" \
-d '{"message": "test", "stream": false}'
```

Expected: Should receive complete JSON response

4. Verify Context Loading

Check the response for:

- `"usedContext": true`
- `"usedCodebase": true`
- `"usedKnowledge": true`
- `"contextError": null`

Troubleshooting

Issue: Still Getting 500 Error

Solution:

1. Check if knowledge base exists: `ls data/ai-knowledge-base.json`
2. If missing, run: `npm run build-knowledge-base`
3. Restart the server

Issue: Slow Response Times

Possible Causes:

1. Context building takes time (expected for first request)
2. Ollama model is slow
3. Large knowledge base

Solutions:

1. Use streaming mode for better UX
2. Disable context for simple queries: `useKnowledge: false, useCodebase: false`
3. Use a faster model: `model: 'phi3:mini'`

Issue: Streaming Not Working

Check:

1. Ensure `stream: true` is set in request
 2. Use `-N` flag with curl for streaming
 3. Check Content-Type header is `text/event-stream`
 4. Verify the application was rebuilt after code changes
-

Future Improvements

Potential Enhancements

1. **Caching:** Cache knowledge base in memory to avoid repeated file reads
2. **Incremental Updates:** Update knowledge base incrementally instead of full rebuild
3. **Compression:** Compress knowledge base file to reduce size
4. **Selective Context:** Only load relevant context based on query type
5. **Progress Indicators:** Add percentage completion for context building

Performance Optimizations

1. **Parallel Processing:** Build context and start Ollama request in parallel
 2. **Context Pruning:** Limit context size based on relevance scores
 3. **Model Selection:** Auto-select faster models for simple queries
 4. **Response Caching:** Cache common queries and responses
-

Related Documentation

- [AI Knowledge System](#) (docs/AI_KNOWLEDGE_SYSTEM.md)
 - [AI Backend Implementation](#) (docs/AI_BACKEND_IMPLEMENTATION_COMPLETE.md)
 - [Codebase Indexing Feature](#) (docs/CODEBASE_INDEXING_FEATURE.md)
-

Pull Request

PR #123: [Enhanced Chat Streaming Support and 500 Error Resolution](#) (<https://github.com/dfulton-thebar/Sports-Bar-TV-Controller/pull/123>)

Branch: `fix/enhanced-chat-streaming-and-knowledge-base`

Base Branch: `optim/ai-perf-security`

Status: Open (Ready for Review)

Conclusion

The 500 Internal Server Error has been successfully resolved by:

1.  Building the missing knowledge base file

2. ☒ Adding streaming support for better UX
3. ☒ Maintaining backward compatibility
4. ☒ Improving error handling and feedback

The chatbot now provides a much better user experience with real-time feedback and successfully loads context from documentation and codebase. Users can see the AI thinking and generating responses in real-time, making the system feel more responsive and transparent.

Next Steps:

1. Review and merge PR #123
2. Deploy to production following the deployment instructions
3. Monitor chatbot performance and user feedback
4. Consider implementing future enhancements listed above