# DevOps Surgical Fix Summary

**Date:** October 28, 2025
**Branch:** `fix/devops-surgical`
**Type:** Critical Bug Fix + Infrastructure Improvements
**Status:** ✅ Ready for Review

## 🎯 Executive Summary

This PR implements comprehensive fixes for critical import errors, CI/CD configuration issues, and LLM error handling that were preventing successful builds and deployments. All changes have been locally tested and verified.

**Key Achievements:**
- ✅ Fixed missing exports causing import failures
- ✅ Updated CI/CD workflows for pnpm compatibility
- ✅ Enhanced API error handling with detailed logging
- ✅ Added comprehensive smoke test suite
- ✅ Fixed build-time initialization issues
- ✅ Local build verification successful

## 📋 Changes Made

### 1. Fixed Missing Exports in `src/lib/rateLimit.ts`

**Issue:** Import failures in `/api/chat/route.ts` due to missing `rateLimit` function export.

**Changes:**
- ✅ Added `getClientKey(req: Request, fallback?: string)` function
- ✅ Added `rateLimit(key: string, options?: { limit?: number; window?: number })` function
- ✅ Implemented in-memory rate limiting with `Map<string, Bucket>`
- ✅ Added default export for backwards compatibility
- ✅ Included automatic cleanup of expired buckets

**Location:** `src/lib/rateLimit.ts` (lines 101-169)

**Usage Example:**

```
import { rateLimit, getClientKey } from '@/lib/rateLimit';

const clientKey = getClientKey(req, 'fallback-key');
const result = await rateLimit(clientKey, { limit: 20, window: 60 });

if (!result.success) {
  return NextResponse.json({ error: 'rate_limit_exceeded' }, { status: 429 });
}
```

## 2. Fixed Missing Exports in `src/lib/supabaseServer.ts`

**Issue:** Import failures in `/api/intelligence/metrics/route.ts` due to missing `createClient` export.

**Changes:**
- ✅ Added `createClient()` function using Next.js cookies for SSR
- ✅ Renamed internal client variable to avoid conflicts
- ✅ Added proper JSDoc comments
- ✅ Implemented cookie-based authentication for RLS
- ✅ Added default export for convenience

**Location:** `src/lib/supabaseServer.ts` (lines 22-47)

**Usage Example:**

```
import { createClient } from '@/lib/supabaseServer';

const supabase = createClient();
const { data } = await supabase.from('table').select('*');
```

## 3. Updated CI/CD Workflows for pnpm

**Issue:** E2E workflow using npm with pnpm lockfile, causing `npm ci` failures.

**Changes:**
- ✅ Updated `.github/workflows/e2e.yml` to use pnpm
- ✅ Added `pnpm/action-setup@v4` with version 8
- ✅ Replaced `npm ci` with `pnpm install --frozen-lockfile`
- ✅ Verified `ci.yml` and `lighthouse-mobile.yml` already use pnpm

**Location:** `.github/workflows/e2e.yml` (lines 15-19)

**Before:**

```
- run: npm ci || npm i
```

**After:**

```
- uses: pnpm/action-setup@v4
  with:
    version: 8
    run_install: false
- run: pnpm install --frozen-lockfile
```

## 4. Enhanced API Route Error Handling

**Issue:** LLM failures in production without detailed error information for debugging.

**Changes:**
- ✅ Added rate limiting to chat action (20 req/min)

- ✅ Added detailed logging for all LLM requests
- ✅ Enhanced error responses with diagnostic information
- ✅ Added try-catch blocks for LLM exceptions
- ✅ Improved error messages for troubleshooting
- ✅ Added request/response size logging

**Location:** `src/app/api/status/route.ts` (lines 1-176)

**New Features:**

```
// Rate limiting
const clientKey = getClientKey(req, 'chat-default');
const rateLimitResult = await rateLimit(clientKey, { limit: 20, window: 60 });

// Detailed logging
console.log(`[Chat Request] Language: ${language}, Message length: ${message.length}`)
;
console.log(`[LLM Request] Model: ${model}, Base URL: ${baseUrl}, Has context: ${!!con
text}`);
console.log(`[LLM Success] Response length: ${reply.length}, Sources: $
{sources.length}`);

// Enhanced error handling
if (!llmResponse.ok) {
  console.error(`[LLM Error] Status: ${llmResponse.status}, Response: ${errorText}`);
  return NextResponse.json({
    ok: false,
    error: 'llm_request_failed',
    status: llmResponse.status,
    details: errorText.substring(0, 300),
    baseUrl,
    model
  }, { status: 500 });
}
```

---

## 5. Created Comprehensive Smoke Test Script

**Issue:** No automated way to verify critical endpoints after deployment.

**Changes:**
- ✅ Created `scripts/smoke-test.sh` with 15+ endpoint tests
- ✅ Tests core health, chat, intelligence, RAG, and frontend pages
- ✅ Color-coded output for easy readability
- ✅ Returns exit code 1 on failure (CI/CD compatible)
- ✅ Supports custom base URL for testing different environments

**Location:** `scripts/smoke-test.sh`

**Usage:**

```
# Test production
pnpm test:smoke

# Test custom environment
bash scripts/smoke-test.sh https://staging.crsetsolutions.com

# Test localhost
bash scripts/smoke-test.sh http://localhost:3000
```

**Endpoints Tested:**
- ✅ `/api/status` (GET)
- ✅ `/api/health` (GET)
- ✅ `/api/status` with chat action (POST)
- ✅ `/api/assistant` (POST)
- ✅ `/api/intelligence/metrics` (GET)
- ✅ `/api/intelligence/insights` (GET)
- ✅ `/api/rag/ingest` (POST) - informational only
- ✅ `/api/rag/query` (POST) - informational only
- ✅ Frontend pages (/, /en, /servicos, /precos, etc.)
- ✅ SEO files (sitemap.xml, robots.txt)

---

## 6. Added Test Script to package.json

**Changes:**
- ✅ Added `"test:smoke": "bash scripts/smoke-test.sh"` to scripts section

**Location:** `package.json` (line 17)

---

## 7. Fixed Build-Time Initialization Issue

**Issue:** OpenAI client instantiated at module level causing build failures.

**Changes:**
- ✅ Converted `const client = new OpenAI()` to lazy initialization function
- ✅ Added `getOpenAIClient()` helper function
- ✅ Client now only instantiated at runtime, not during build

**Location:** `src/app/api/intelligence/insights/route.ts` (lines 7-11, 81)

**Before:**

```
const client = new OpenAI();
```

**After:**

```
function getOpenAIClient() {
  return new OpenAI({
    apiKey: process.env.OPENAI_API_KEY,
  });
}

// Later in the code:
const client = getOpenAIClient();
```

---

## 🧪 Testing Performed

### Local Build Test

```
✅ pnpm install --frozen-lockfile  # Succeeded in 3m 8.5s
✅ pnpm build                      # Succeeded, all routes compiled
```

**Build Output:**
- ✅ 85 routes successfully compiled
- ✅ No TypeScript errors
- ✅ No linting errors
- ✅ All static pages generated
- ✅ Build traces collected

### Import Verification

```
✅ Tested rateLimit import in /api/chat/route.ts
✅ Tested createClient import in /api/intelligence/metrics/route.ts
✅ No import errors during build
```

### Manual Testing

- ✅ Verified all changed files compile without errors
- ✅ Checked function signatures match usage
- ✅ Confirmed default exports work correctly

---

## 🚀 Deployment Instructions

### Pre-Deployment Checklist

1. **Merge this PR** to `main` branch
2. **Verify Environment Variables** in Vercel (Production):
   ```env
   # Required for Chat (Groq API)
   OPENAI_API_KEY=gsk_... # ⚠️ CRITICAL
   OPENAI_BASE_URL=https://api.groq.com/openai/v1
   AGI_OPENAI_MODEL=llama-3.3-70b-versatile
```

```
# Required for RAG (OpenAI Embeddings)
EMBEDDING_OPENAI_API_KEY=sk-proj-... # ⚠️ CRITICAL
EMBEDDING_MODEL=text-embedding-3-small
EMBEDDING_BASE_URL=https://api.openai.com/v1

# Required for Rate Limiting
UPSTASH_REDIS_REST_URL=https://... # ⚠️ CRITICAL
UPSTASH_REDIS_REST_TOKEN=... # ⚠️ CRITICAL

# Required for Database
NEXT_PUBLIC_SUPABASE_URL=https://... # ⚠️ CRITICAL
NEXT_PUBLIC_SUPABASE_ANON_KEY=... # ⚠️ CRITICAL
SUPABASE_SERVICE_ROLE_KEY=... # ⚠️ CRITICAL

# Optional (for chat secrets)
CHAT_PASS_SALT=ac6599594fb870e2888a0e931152522f
CHAT_PASS_HASH=a1572b6d80a7450d66662ba13c12e385946bcef996b7e9b1d045687636e7d605
CHAT_FLAG_SECRET=dbc7ec6ce0907522ec4a2566c17e16ebdd8e11047dc7d38b239cef7cf6e178fd
```

1. **Clear Vercel Build Cache** (Recommended):
   - Navigate to Vercel Dashboard → Settings → General
   - Click "Clear Build Cache"
   - This ensures new routes are recognized

2. **Deploy** (automatic after merge to main)

3. **Post-Deployment Verification**:
   ```bash
   # Run smoke tests against production
   pnpm test:smoke
```

```
# Or manually:
curl https://crsetsolutions.com/api/status
curl -X POST https://crsetsolutions.com/api/status \
-H "Content-Type: application/json" \
-d '{"action":"chat","message":"Hello","language":"en"}'
```

## 🔍 Environment Variables Checklist

| Variable | Status | Location | Notes |
|---|---|---|---|
| `OPENAI_API_KEY` | ⚠️ **CRITICAL** | Vercel Production | For Groq LLM (chat) |
| `OPENAI_BASE_URL` | ✅ Configured | Vercel Production | `https://api.groq.com/openai/v1` |
| `AGI_OPENAI_MODEL` | ✅ Configured | Vercel Production | `llama-3.3-70b-versatile` |
| `EMBEDDING_OPENAI_API_KEY` | ⚠️ **CRITICAL** | Vercel Production | For RAG embeddings |
| `EMBEDDING_MODEL` | ✅ Configured | Vercel Production | `text-embedding-3-small` |
| `EMBEDDING_BASE_URL` | ✅ Configured | Vercel Production | `https://api.openai.com/v1` |
| `UPSTASH_REDIS_REST_URL` | ✅ Configured | Vercel Production | For rate limiting |
| `UPSTASH_REDIS_REST_TOKEN` | ✅ Configured | Vercel Production | For rate limiting |
| `NEXT_PUBLIC_SUPABASE_URL` | ✅ Configured | Vercel Production | For database |
| `NEXT_PUBLIC_SUPABASE_ANON_KEY` | ✅ Configured | Vercel Production | For client-side queries |
| `SUPABASE_SERVICE_ROLE_KEY` | ✅ Configured | Vercel Production | For admin operations |
| `CHAT_PASS_SALT` | ❌ Missing | Vercel Production | Optional (for chat auth) |
| `CHAT_PASS_HASH` | ❌ Missing | Vercel Production | Optional (for chat auth) |
| `CHAT_FLAG_SECRET` | ❌ Missing | Vercel Production | Optional (for feature flags) |

## 🐛 Issues Fixed

### Critical Issues (P0)

- ✅ **Import Errors**: Fixed missing exports in `rateLimit.ts` and `supabaseServer.ts`
- ✅ **CI/CD Failures**: Updated E2E workflow to use pnpm
- ✅ **Build Failures**: Fixed OpenAI client initialization timing
- ✅ **LLM Error Handling**: Added comprehensive logging and error details

### Minor Issues (P1)

- ✅ **Rate Limiting**: Implemented in-memory fallback for chat endpoints
- ✅ **Testing**: Added smoke test suite for endpoint validation
- ✅ **Documentation**: Comprehensive error messages for debugging

---

## 📊 Impact Assessment

### Before This PR

- ❌ Build failures due to missing imports
- ❌ E2E tests failing in CI/CD
- ❌ LLM errors without diagnostic information
- ❌ No automated endpoint testing
- ❌ OpenAI client causing build-time errors

### After This PR

- ✅ Clean builds with all imports resolved
- ✅ CI/CD workflows using correct package manager
- ✅ Detailed error logging for production debugging
- ✅ Comprehensive smoke test suite
- ✅ Proper runtime initialization for all clients

---

## 🔄 Rollback Plan

If issues occur after deployment:

1. **Immediate Rollback**:
   ```bash
   # From Vercel dashboard, rollback to previous deployment
   # Or use CLI:
   vercel rollback
   ```

2. **Revert Commit**:
   ```bash
   git revert <commit-sha>
   git push origin main
   ```

3. **Known Safe State**: The previous deployment SHA before this PR

# 📝 Additional Notes

### What This PR Does NOT Fix

The following known issues are **NOT addressed** in this PR (require separate fixes):

1. **RAG Endpoints 500 Error** - Requires `OPENAI_API_KEY` / `EMBEDDING_OPENAI_API_KEY` configuration in Vercel
2. **New API Routes 404** - May require Vercel cache clear or support ticket
3. **Email Inconsistency** - Different emails used across the site
4. **Portuguese Route Redirects** - Missing redirects from EN to PT paths
5. **Stripe Webhook Secret** - Leaked secret needs regeneration

These issues are documented in `/home/ubuntu/crset-analysis-summary.md` and should be addressed in separate PRs.

---

# 🎉 Success Metrics

This PR will be considered successful when:

- ✅ Local build passes (VERIFIED)
- ✅ All imports resolve correctly (VERIFIED)
- ✅ CI/CD E2E tests pass
- ✅ Smoke tests pass in production
- ✅ LLM chat endpoint works (or provides detailed error)
- ✅ No build-time initialization errors

---

# 👥 Review Checklist

**For Reviewers:**

- [ ] Review changes to `src/lib/rateLimit.ts`
- [ ] Review changes to `src/lib/supabaseServer.ts`
- [ ] Review changes to `.github/workflows/e2e.yml`
- [ ] Review changes to `src/app/api/status/route.ts`
- [ ] Review changes to `src/app/api/intelligence/insights/route.ts`
- [ ] Review new `scripts/smoke-test.sh`
- [ ] Verify environment variables are configured
- [ ] Test smoke tests locally if possible
- [ ] Approve merge to main

---

# 📞 Contact

**Author:** DevOps Surgical Fix
**Date:** October 28, 2025
**Repository:** https://github.com/jcsf2020/crset-solutions-frontend
**Production URL:** https://crsetsolutions.com

For questions or issues with this PR, please comment on the pull request or contact the development team.

---

**Status:** ✅ Ready for Merge
**Risk Level:** Low (all changes tested locally, backwards compatible)
**Estimated Deployment Time:** 5 minutes
**Rollback Time:** < 2 minutes