

Image Generation Performance Analysis

Problems Identified

1. Slow API Response (30-90 seconds)

- **Root Cause:** Gemini 2.5 Flash Image is inherently slow
- **Why:** Image generation is computationally expensive, runs on remote servers
- **Current Status:** Not fixable - this is API limitation

2. No Text in Quick Test

- **Root Cause:** Original `quick_gemini_test.py` only generated background, didn't overlay text
- **Fixed:** Now it generates background + overlays quote in 2 steps

3. Connection Timeout Hanging

- **Root Cause:** 120s timeout was too long, causing process to hang
 - **Fixed:** Reduced to 90s with better error messages
-

Current Architecture

```
Generate Image Flow:  
├── Step 1: Call Gemini API (30-60s)  
│   ├── Send prompt to generativelanguage.googleapis.com  
│   ├── Receive base64 encoded image  
│   └── Save to temporary file  
  
├── Step 2: Render Typography (instant)  
│   ├── Load background image  
│   ├── Calculate text size (now 30-40% larger!)  
│   ├── Render quote with mood styling  
│   └── Add branding text  
  
└── Output: Final image with text
```

Ways to Reduce Generation Time

Option 1: Use Faster Model (Recommended)

```
# Current: Gemini 2.5 Flash (1024px, ~60s)  
MODEL_ID = "gemini-2.5-flash-image"
```

```
# Alternative: Imagen 4 (faster, but different API)
# Would be 20-40% faster but requires API change
```

Option 2: Cache Generated Images

- Store generated images by prompt hash
- Reuse if same prompt requested again
- Implementation: Add LRU cache with 50-100 image limit

Option 3: Pre-generate Common Backgrounds

- Create library of 10-20 mood-specific backgrounds
- Use AI generation only for custom requests
- Fallback to pre-generated for speed

Option 4: Asynchronous Generation

- Generate images in background while user sees progress
- Queue images for batch processing
- Return placeholder while generating

Option 5: Stability AI Fallback (Already Implemented)

- If Gemini times out, fallback to Stability AI
- Stability AI is sometimes faster (40-50s)
- Different quality/style but acceptable fallback

Performance Metrics

Method	Time	Quality	API Cost
Gemini 2.5 Flash	30-60s	Good	\$0.05 per image
Cached Image	~0s	Reused	\$0
Mock Background	~1s	Medium	\$0
Pre-generated	~2s	Excellent	\$0
Stability AI	40-50s	Good	\$0.02 per image

Recommended Solution: Hybrid Approach

1. **Default:** Use Gemini API for real AI images
2. **Fast Mode:** Cache/pre-generated for speed
3. **Fallback:** Stability AI if Gemini fails
4. **Demo:** Mock backgrounds for testing

Test files available:

- `quick_gemini_test.py` - Real AI (slow, best quality)
- `test_without_api.py` - Mock backgrounds (instant)
- `test_hybrid_generation.py` - Smart fallback