

Qwen2.5-VL Not Loading in LM Studio? Fix mmproj and Vision Errors

February 26, 2026 · by Mark Bartlett

[Download this guide as PDF](#)

Quick Answer: The most common reason Qwen2.5-VL doesn't work in LM Studio: you downloaded the model GGUF but not the mmproj file. Vision models need both files in the same folder. Download mmproj-model-f16.gguf from the lmstudio-community repo, place it next to your model GGUF, restart LM Studio, and look for the yellow eye icon. If you see 'model type qwen2_5_vl not supported', update LM Studio to 0.3.10+. If the model crashes on image input, switch from a BF16 mmproj to FP16, and avoid iq4_xs quantizations from third-party repos.

 **Related:** [Qwen2.5-VL Setup Guide \(happy path\)](#) · [Vision Models Locally](#) · [Qwen Models Guide](#) · [VRAM Requirements](#) · [LM Studio Tips & Tricks](#)

We have a [full setup guide for Qwen2.5-VL in LM Studio](#). This article is for when that didn't work. You followed the steps, the model loaded, and either vision isn't available or something crashed.

Every error below is documented from LM Studio's bug tracker and HuggingFace discussions. These aren't hypothetical – they're the issues people actually hit.

Error 1: No eye icon, no image button

What you see: The model loads fine for text chat. But there's no yellow eye icon next to the model name, no image button in the chat input, and no way to attach a photo.

Cause: The mmproj file is missing or not in the right folder.

Qwen2.5-VL is a two-file model. The main GGUF holds the language model. A separate file called `mmproj-model-f16.gguf` holds the vision encoder projector – the piece that translates images into tokens the language model can process. Without it, LM Studio loads the model as text-only.

Fix

1. Go to the [lmstudio-community Qwen2.5-VL repo](#) on HuggingFace
2. Download `mmproj-model-f16.gguf` (about 1.35 GB)
3. Place it in the **same directory** as your model GGUF

The default LM Studio model directory is:

OS	Path
macOS	<code>~/lmstudio/models/lmstudio-community/Qwen2.5-VL-7B-Instruct-GGUF/</code>
Windows	<code>C:\Users\<>you>\.lmstudio\models\lmstudio-community\Qwen2.5-VL-7B-Instruct-GGUF\</code>
Linux	<code>~/lmstudio/models/lmstudio-community/Qwen2.5-VL-7B-Instruct-GGUF/</code>

Your folder should contain both files:

```
Qwen2.5-VL-7B-Instruct-GGUF/
├─ Qwen2.5-VL-7B-Instruct-Q4_K_M.gguf    (4.68 GB)
└─ mmproj-model-f16.gguf                (1.35 GB)
```

1. Restart LM Studio (or unload and reload the model)
2. Look for the yellow eye icon next to the model name in the chat dropdown

If the eye icon appears, vision is active. If it doesn't, check the next section.

Still no eye icon after adding the mmproj?

Three things to check:

Wrong directory. The mmproj file must be in the exact same folder as the model GGUF. Not a parent folder, not a sibling folder. Same folder.

Wrong filename. LM Studio detects mmproj files by filename pattern. The file must contain "mmproj" in its name. If you renamed it or downloaded a variant with an unusual name (like `vision-encoder.gguf`), LM Studio won't detect it.

Non-ASCII characters in the path. This is a [documented bug](#). If your model path contains Chinese characters, accented letters, or other non-ASCII text, the mmproj file fails to load silently. Move the model to a path with only ASCII characters.

Error 2: "Model type qwen2_5_vl not supported"

What you see: LM Studio throws `ValueError: Model type qwen2_5_vl not supported` when you try to load the model.

Cause: Your LM Studio version is too old. Qwen2.5-VL support was added after the initial Qwen2-VL support:

Platform	Backend	Minimum version
macOS	MLX engine	0.3.10+ (MLX runtime 0.4.0+)
Windows	llama.cpp CUDA	0.3.14+ (runtime v1.25+)
Linux	llama.cpp	0.3.14+ (runtime v1.25+)

Fix

1. Update LM Studio: Help → Check for Updates, or download the latest from lmstudio.ai
2. On Mac, make sure you're using the **MLX backend**, not llama.cpp. The MLX engine added Qwen2.5-VL support in version 0.4.0 (merged Feb 2025). Check Settings → Runtime to confirm.
3. If you're on an older macOS version that can't run the latest LM Studio, Qwen2-VL (not 2.5) may work as a fallback with fewer capabilities.

Note: Qwen2-VL (the older model) works on earlier LM Studio versions. Qwen2.5-VL requires the newer backends. These are different architectures – LM Studio treats them as separate model types.

Error 3: Model crashes on load (exit code 1844674407...)

What you see: LM Studio starts loading the model, then crashes with a long exit code number, sometimes `18446744072635810000` or similar.

Cause: Almost always a bad quantization. The most common culprit is Mungert's `iq4_xs` quantization of Qwen2.5-VL, which has [known loading issues](#) in LM Studio.

Fix

Delete the problematic GGUF and download one from the **lmstudio-community** repos instead:

Size	Repo
3B	<code>lmstudio-community/Qwen2.5-VL-3B-Instruct-GGUF</code>
7B	<code>lmstudio-community/Qwen2.5-VL-7B-Instruct-GGUF</code>

Size	Repo
32B	<code>lmstudio-community/Qwen2.5-VL-32B-Instruct-GGUF</code>
72B	<code>lmstudio-community/Qwen2.5-VL-72B-Instruct-GGUF</code>

These are quantized by bartowski and tested against LM Studio's backends. Stick with Q4_K_M for the safest balance of size and compatibility.

Third-party quantizations (especially `iq4_xs`, `iq3_xxs`, and abilitated variants) may work but have a higher failure rate. If you want an abilitated variant, use one that includes its own mmproj files – check the repo's Files tab before downloading.

Error 4: Model crashes on image input (not on text)

What you see: The model loads fine. Text chat works. But when you attach an image and send a message, LM Studio crashes or hangs.

Cause: Usually one of three things:

BF16 mmproj file

Some repos provide mmproj files in BF16 (bfloat16) format. The CUDA backend in llama.cpp [does not support BF16 for the im2col operation](#) used in vision processing. The model loads because the language model is fine, but it crashes when the vision encoder tries to process an image.

Fix: Download the FP16 (float16) version of the mmproj file instead. Look for `mmproj-model-f16.gguf` rather than `mmproj-model-bf16.gguf`. The lmstudio-community repos provide FP16 by default.

VRAM overflow on image processing

Vision inference spikes VRAM temporarily. The vision encoder needs activation memory on top of the model weights and KV cache. High-resolution images consume more.

Fix:

- Reduce context length to 2048-4096 in model settings
- Drop to a smaller quantization (Q6_K → Q4_K_M)
- Resize large images before sending (under 1024px on the longest side)

- Check Settings → Chat → Image Inputs for the max dimension setting

Memory requirements for vision

The mmproj file adds ~1.35 GB on top of the language model:

Model	Model GGUF (Q4_K_M)	mmproj (FP16)	Total VRAM needed
3B	1.93 GB	1.34 GB	~5-6 GB
7B	4.68 GB	1.35 GB	~8-10 GB
32B	19.9 GB	1.38 GB	~23-24 GB
72B	47.4 GB	1.41 GB	~50-52 GB

If your GPU is within 1-2 GB of these numbers, image processing will push you over. Drop to a lower quantization or a smaller model.

Error 5: “I cannot view images” or model ignores the image

What you see: No crash. The model responds, but it says something like “I cannot view images” or “I don’t see any image in our conversation” – even though you attached one.

Cause: The mmproj file isn’t actually connected, even if it’s present. This happens when:

1. **The mmproj loaded for a different model.** If you have multiple vision models, LM Studio may associate the wrong mmproj with the wrong model. Unload all models and reload just the Qwen2.5-VL model.
2. **You loaded the model before placing the mmproj file.** LM Studio checks for the mmproj during model load, not continuously. If you added the file while the model was already running, unload and reload.
3. **LM Studio is using the wrong backend.** On Mac, the MLX backend handles Qwen2.5-VL vision correctly. The llama.cpp backend may load the language model but fail to connect the vision encoder. Switch to MLX in Settings → Runtime.

Fix

1. Unload the model completely
2. Confirm the mmproj file is in the same folder as the model GGUF

3. Reload the model
4. Verify the yellow eye icon appears
5. Test with a simple prompt: "What is in this image?" with a clear photo

If the model still ignores images, check the LM Studio developer console (View → Toggle Developer Tools) for vision-related error messages.

Error 6: Vision works but output is garbage

What you see: The model acknowledges the image but produces incoherent, repetitive, or completely wrong descriptions.

Cause: Several possibilities:

Context length too short

Each image gets converted to visual tokens. A standard image produces 500-1,500 tokens. If your context length is set to 2048 and the image eats 1,200 tokens, the model has almost no room to generate a response.

Fix: Set context length to at least 4096 for vision use. 8192 is safer for detailed image analysis. Check Model Settings → Context Length.

Wrong chat template

If LM Studio applies the wrong chat template, the model gets confused prompts. This mostly happens with third-party quantizations that don't include the correct metadata.

Fix: Use models from the `lmstudio-community` repos, which include correct template metadata. If you must use a third-party GGUF, check Model Settings → Chat Template and select "Qwen2VL" or let LM Studio auto-detect.

Heavily quantized model

At Q2 or IQ2 quantization, vision quality degrades noticeably. The model may recognize objects but produce garbled descriptions or miss details.

Fix: Use Q4_K_M or higher for vision tasks. The quality jump from Q3 to Q4 is significant for multimodal models.

Error 7: VRAM not released after unloading

What you see: After unloading a Qwen2.5-VL model (or after a crash), VRAM stays consumed. Loading a new model fails because memory is still held.

Cause: Known bug in LM Studio, primarily on Linux. The vision encoder's memory allocation isn't always cleaned up on model unload.

Fix

1. Restart LM Studio completely (not just unload the model)
 2. If that doesn't free VRAM, check for orphaned processes: `ps aux | grep lmstudio` (Linux/Mac) or Task Manager (Windows)
 3. As a last resort, reboot
-

Quick-reference checklist

Before reporting a bug, run through this:

Check	What to do
LM Studio version	0.3.14+ on Windows/Linux, 0.3.10+ on Mac
mmproj file present	<code>mmproj-model-f16.gguf</code> in same folder as model GGUF
Eye icon visible	Yellow eye next to model name in chat dropdown
Quantization source	lmstudio-community repo (avoid iq4_xs from third parties)
mmproj format	FP16, not BF16 (especially on Windows/NVIDIA)
Context length	4096+ for vision use
Path characters	ASCII only in model directory path
Backend (Mac)	MLX engine, not llama.cpp

Which mmproj to download

Two options exist for most repos:

mmproj type	File size (7B)	When to use
FP16 (<code>mmproj-model-f16.gguf</code>)	~1.35 GB	Default choice. Works everywhere.
Q8_0 (<code>mmproj-Q8_0.gguf</code>)	~0.7 GB	Saves ~600 MB. Slightly lower vision quality. No imatrix support.

FP16 is the safe default. The mmproj file is small relative to the model – saving 600 MB by quantizing it rarely makes a meaningful difference. Use Q8_0 only if you're within 1 GB of your VRAM limit and need every byte.

Don't use BF16 mmproj files on NVIDIA GPUs. The CUDA backend [crashes during image processing](#) with BF16 mmproj. This is a llama.cpp limitation, not an LM Studio bug. FP16 works identically in practice.

When nothing works

If you've tried everything above and Qwen2.5-VL still won't work:

- 1. Try Ollama instead.** `ollama pull qwen2.5vl:7b` handles the mmproj automatically. No file management needed. See our [Vision Models guide](#) for Ollama setup.
- 2. Try a different vision model.** Gemma 3 4B and 12B work reliably in LM Studio with fewer quirks. LLaVA 1.6 is the oldest and most stable option.
- 3. Report the bug.** File an issue at [lmstudio-ai/lmstudio-bug-tracker](#) with your LM Studio version, OS, GPU, model file names, and the exact error message. Include a screenshot of your model folder showing both files.

Bottom line

90% of Qwen2.5-VL problems in LM Studio come down to the mmproj file. It's either missing, in the wrong folder, or the wrong format. Download `mmproj-model-f16.gguf` from the [lmstudio-community repo](#), drop it next to your model GGUF, restart, and look for the eye icon.

The other 10% is version issues (update LM Studio) or bad quantizations (use lmstudio-community repos, not third-party ones). If the model loads and the eye icon appears, vision works.

Get notified when we publish new guides.

[Subscribe – free, no spam](#)

Source: <https://insiderllm.com/guides/qwen25-vl-lm-studio-troubleshooting/>

Free guides for running AI locally