# Intel Arc GPUs for Local AI: The Underdog Option That Actually Works

February 24, 2026 · by Mark Bartlett

[Download this guide as PDF](#)

> **Quick Answer:** The Intel Arc A770 16GB is the cheapest way to get 16GB of VRAM for local AI, around $250 used. Software support through IPEX-LLM (Ollama integration), llama.cpp's SYCL backend, and OpenVINO for image generation is real and improving. Expect 40-55 tok/s on 7B models with llama.cpp SYCL, which is slower than an RTX 3060 12GB on CUDA but fast enough for interactive use. The A750 8GB isn't worth it: at 8GB VRAM, just buy an RTX 3060 12GB with better software support. Buy the A770 if you want VRAM on a budget and can tolerate some software friction.

📚 **More on this topic:** [GPU Buying Guide](#) · [AMD vs NVIDIA for Local AI](#) · [What Can You Run on 16GB VRAM](#)

Nobody talks about Intel Arc for local AI. When people ask "which GPU should I buy for running LLMs," the answer is always NVIDIA first, AMD second, Intel never.

That's mostly fair. NVIDIA's CUDA ecosystem is dominant. AMD's ROCm has caught up enough to be viable. Intel's software stack is the youngest of the three, with the smallest community and the most rough edges.

But Intel has one card that keeps showing up in budget discussions: the Arc A770 16GB. Sixteen gigabytes of VRAM for around $250 used. No NVIDIA card at that price comes close — the cheapest 16GB option from NVIDIA is the RTX 4060 Ti 16GB at $450+, and the popular RTX 3060 only has 12GB.

So the question isn't whether Intel Arc is better than NVIDIA. It isn't. The question is whether 16GB of VRAM at $250 is worth the software headaches.

---

## The Hardware: Two Cards Worth Discussing

Intel's Arc A-series (Alchemist) launched in late 2022 to mixed reviews. Gaming performance was inconsistent. Drivers were rough. But two cards landed in an interesting spot for AI work:

## Arc A770 16GB — The One That Matters

| Spec | Arc A770 16GB |
| --- | --- |
| VRAM | 16GB GDDR6 |
| Memory Bandwidth | 560 GB/s |
| Xe Cores | 32 |
| TDP | 225W |
| New Price | ~$280-350 |
| Used Price | ~$250-265 |
| Original MSRP | $329 (2022) |

The A770 is the entire reason this article exists. 16GB of VRAM at ~$250 used is a price point that doesn't exist anywhere else in the GPU market. For context:

| GPU | VRAM | Used Price | $/GB |
| --- | --- | --- | --- |
| **Arc A770** | **16GB** | **~$255** | **$16/GB** |
| RTX 3060 | 12GB | ~$180 | $15/GB |
| RTX 3060 Ti | 8GB | ~$200 | $25/GB |
| RTX 4060 | 8GB | ~$260 | $32/GB |
| RX 6800 XT | 16GB | ~$300 | $19/GB |

The A770 wins on raw VRAM capacity at its price. That 16GB lets you run 14B parameter models at Q4 quantization — models like Qwen 2.5 14B, Llama 3.1 8B at higher quants, or even squeeze in a 30B model at aggressive Q2 compression. An 8GB card cuts you off at 7B models.

## Arc A750 8GB — Skip It

| Spec | Arc A750 8GB |
| --- | --- |
| VRAM | 8GB GDDR6 |
| Memory Bandwidth | 512 GB/s |
| Xe Cores | 28 |
| TDP | 225W |

| Spec | Arc A750 8GB |
|---|---|
| New Price | ~$160-200 |
| Used Price | ~$146 |

At 8GB, the A750 loses its only advantage. An RTX 3060 12GB costs ~$180 used, has 12GB of VRAM instead of 8GB, and runs on CUDA with flawless software support. The A750 saves you $30-40 but gives you less VRAM and worse compatibility. Not worth it.

The newer **Arc B580 12GB** ($249 new, Xe2 architecture) is a better buy than the A750: more VRAM, newer architecture, better driver support. If you want an Intel card under $250, get the B580 instead. This article focuses on the A770 16GB, which remains the best deal for maximum VRAM on a budget.

## The Software Story: Three Paths to Inference

Running LLMs on Intel Arc requires different software than NVIDIA's plug-and-play CUDA stack. There are three paths that actually work.

### Path 1: IPEX-LLM + Ollama (Easiest)

[IPEX-LLM](#) (formerly BigDL-LLM) is Intel's optimized LLM inference library. It wraps llama.cpp with Intel-specific optimizations and provides a modified Ollama binary that works directly with Arc GPUs.

**Current status:** The original Intel repo was archived in January 2026, but development continues under a community fork at `ipex-llm/ipex-llm`. The latest stable release is v2.2.0 (April 2025) with nightly v2.3.0b builds adding Ollama v0.9.3 support.

**Setup (Portable Zip — Windows):**

```
1. Update Intel GPU drivers to latest
2. Download IPEX-LLM Ollama portable zip from releases
3. Extract to any folder
4. Run start-ollama.bat
5. ollama run llama3.1:8b
```

That's it. No oneAPI installation, no environment variables, no building from source. The portable zip bundles everything.

**Setup (Linux):**

```
# Ensure kernel 6.8+ and latest Intel GPU drivers
# Download portable .tgz from ipex-llm releases
tar xzf ipex-llm-ollama-portable-*.tgz
cd ipex-llm-ollama
./start-ollama.sh

# In another terminal:
ollama run deepseek-r1:7b
```

**What it supports:** 70+ model architectures including Llama, Mistral, Qwen, DeepSeek, Phi, Gemma, and ChatGLM. FP4/INT4/FP8 quantization. Even some multimodal models like Qwen-VL in recent builds.

**The catch:** You're running a modified Ollama, not the official release. Updates lag behind mainline Ollama. And since Intel archived the original repository, long-term maintenance depends on the community fork continuing.

## Path 2: llama.cpp SYCL Backend (Most Stable Long-Term)

The SYCL backend is merged into mainline llama.cpp. Not a fork, not a patch — it's part of the official codebase, which makes it the safest bet long-term.

**Current status:** Actively maintained. Verified with Intel oneAPI 2025.2.1, 2025.1, and 2024.1. Supports all Arc discrete GPUs and integrated GPUs from 11th Gen onward.

**Building from source (Linux):**

```
# Install Intel oneAPI Base Toolkit first
source /opt/intel/oneapi/setvars.sh

cmake -B build -DGGML_SYCL=ON \
  -DCMAKE_C_COMPILER=icx \
  -DCMAKE_CXX_COMPILER=icpx
cmake --build build --config Release -j $(nproc)

# Run inference
./build/bin/llama-cli -m model.gguf -ngl 99 -p "Hello"
```

**Building from source (Windows):**

```
cmake -B build -G "Ninja" -DGGML_SYCL=ON ^
  -DCMAKE_C_COMPILER=cl ^
  -DCMAKE_CXX_COMPILER=icx ^
  -DCMAKE_BUILD_TYPE=Release
cmake --build build --config Release
```

Requires the Intel oneAPI Base Toolkit (or the smaller Intel Deep Learning Essentials package).

**Performance note:** Legacy quantization formats (Q4_0, Q4_1, Q5_0, Q5_1, Q8_0) perform best on the SYCL backend. K-quants (Q4_K_M, Q5_K_M), which are what most people actually use, are not as well optimized yet. This matters: the community standard is K-quants, and Intel's backend runs them slower than it should.

### Path 3: Ollama Native SYCL (Newest, Least Tested)

Ollama v0.17 (February 2026) added improved Intel Arc support through SYCL integration. Multiple pull requests have landed to support Intel GPUs natively, without needing IPEX-LLM.

**Current reality:** It works, but it's described as "less mature than NVIDIA or AMD" in early coverage. If you want the most reliable experience today, use IPEX-LLM's Ollama portable zip. If you want to bet on the future, native Ollama SYCL support is the path that will eventually win.

---

# Benchmarks: How Fast Is It Actually?

Intel's official numbers, independent benchmarks, and user reports tell three different stories. I'll show all of them.

### LLM Inference (tok/s)

**llama.cpp SYCL backend on Arc A770:**

| Model | Quant | Arc A770 (SYCL) | RTX 3060 12GB (CUDA) | RTX 4060 8GB (CUDA) |
|---|---|---|---|---|
| Llama 2 7B | Q4_0 | ~55 tok/s | ~65 tok/s | ~70 tok/s |
| 7B models (general) | Q4_K_M | ~40-45 tok/s | ~55-60 tok/s | ~60-65 tok/s |
| 14B models | Q4_K_M | ~18-22 tok/s | ~25-30 tok/s | Can't fit (8GB) |

**IPEX-LLM (Intel's optimized path):**

Intel's official benchmark claims ~70 tok/s on Mistral 7B (INT4, FP16 arithmetic) on the A770. Independent user testing on a real A770 with an i3-12100 and Ubuntu shows:

| Model Size | IPEX-LLM (user-reported) |
|---|---|
| 1-3B | ~25-30 tok/s |
| 7-8B | ~14-16 tok/s |
| 12-14B | ~10 tok/s |

The gap between Intel's marketing numbers and real-world results is substantial. CPU bottlenecking (the i3-12100 is weak), driver versions, and specific quantization methods all contribute. A better CPU would close part of that gap, but the A770 is still slower than NVIDIA at the same model size.

## Where the A770 Wins

**It doesn't win on speed.** At 7B models, both the RTX 3060 and RTX 4060 are faster. The 4060 is roughly 30-40% faster at the same model size.

**It wins on what you can run at all.** The RTX 4060 has 8GB. The RTX 3060 has 12GB. The A770 has 16GB. When you want to run a 14B model (Qwen 2.5 14B, DeepSeek R1 14B, Llama 3.1 8B at Q8), the A770 loads it while the 4060 can't and the 3060 struggles at higher quants.

| Scenario | Arc A770 16GB | RTX 3060 12GB | RTX 4060 8GB |
|---|---|---|---|
| 7B Q4_K_M | Runs (slower) | Runs (faster) | Runs (fastest) |
| 14B Q4_K_M | Runs (~20 tok/s) | Tight fit, might work | Won't fit |
| 14B Q5_K_M | Runs | Won't fit | Won't fit |
| 30B Q2_K | Barely fits | Won't fit | Won't fit |
| Price (used) | ~$255 | ~$180 | ~$260 |

**20 tok/s on a 14B model is usable for interactive chat.** 0 tok/s because the model doesn't fit is not.

## Image Generation

For Stable Diffusion and Flux, the A770 works through OpenVINO but is significantly slower than NVIDIA:

| GPU | SDXL (it/s) | SD 1.5 (images/min) |
| --- | --- | --- |
| Arc A770 | ~3.4 | ~15.4 |
| RTX 3060 12GB | ~6.6 | ~30+ |
| RTX 4060 8GB | ~8.1 | ~35+ |

The A770 is roughly half the speed of an RTX 3060 for SDXL. If image generation is your primary use case, buy NVIDIA. The A770's image gen performance works, but you won't enjoy it.

## What Actually Works Today

### Working well

**Ollama via IPEX-LLM** is the easiest path. The portable zip approach gives you a working Ollama install in about two minutes: download, extract, run. Models load, inference works, and the Ollama API is compatible with tools that expect Ollama.

**llama.cpp SYCL** is merged into mainline and actively maintained. It produces correct output and delivers usable speed. Legacy quants (Q4_0, Q8_0) run faster than K-quants on this backend.

**OpenVINO** handles image generation. Stable Diffusion 1.5, SDXL, and Flux all work through it. ComfyUI has an official Intel OpenVINO node, and Automatic1111 WebUI supports the OpenVINO backend too.

**PyTorch 2.5+** has native Intel XPU support. Basic inference and some training workloads run on Arc GPUs through the `torch.xpu` device.

### Working, but rough

**Ollama native SYCL** landed in v0.17 but it's early. Expect occasional issues that IPEX-LLM's more tested integration handles better.

**K-quant performance** is a real problem. The SYCL backend handles K-quants (Q4_K_M, Q5_K_M) worse than legacy quants. Since K-quants are the community standard, real-world performance often lags benchmarks that use Q4_0.

**Dual GPU setups** don't work well. Running two A770s together through IPEX-LLM actually slows inference down. Stick with a single card.

**Linux drivers** remain unstable. Arc GPUs require kernel 6.8+. The ongoing `i915` to `xe` driver transition causes breakage on some systems, and if your iGPU and dGPU are both active, SYCL may fail to initialize. You might need to disable the iGPU in BIOS.

### Doesn't work (yet)

**Most training frameworks** are CUDA-only. LoRA fine-tuning tools like Unsloth, DreamBooth training scripts, and many HuggingFace training examples won't run on Arc. Inference works; training is off the table.

**Video generation** tools (AnimateDiff, Stable Video Diffusion) require CUDA. No workaround.

**Hardcoded `torch.cuda` projects** won't detect your Arc GPU. If a project calls `torch.cuda.is_available()` without an XPU abstraction layer, it just skips your hardware. This is common in older projects and research code.

`torch.compile` causes segfaults on Arc after PyTorch 2.8.0. Known bug, not fixed yet.

**ZLUDA** is dead. Both AMD and Intel dropped support for this CUDA translation layer. It would have let CUDA code run on non-NVIDIA GPUs, but nobody wanted to maintain it.

---

# The Driver Situation

Intel Arc drivers have improved since the rough 2022 launch, but they still trail NVIDIA in stability.

**On Windows:** Drivers are fine now. Regular updates have fixed the worst gaming issues, and for AI workloads, the IPEX-LLM portable zip bundles what it needs. Driver version matters less for AI than for gaming.

**On Linux:** This is where pain lives.

- Requires kernel 6.8+ for Alchemist (A-series) GPUs
- Ubuntu LTS ships older kernels by default, so you may need to upgrade
- The `i915` to `xe` driver transition is mid-stream and causes breakage
- Performance is 25-50% lower than Windows in some GPU workloads
- No `nvidia-smi` equivalent; Intel's monitoring tools are limited
- System freezes have been reported under heavy GPU load
- Arc A770/A750 were broken on Linux 6.19 Git kernel (no display output), since fixed

**Practical advice:** If you're buying an A770 for AI on Linux, use Ubuntu 24.04 with kernel 6.8+ and test everything before committing to the setup. Keep a known-working kernel as a fallback boot option.

## Buy the A770 if…

You need 16GB VRAM and your GPU budget is under $300. The A770 is the only card at that price. Next cheapest 16GB option is the RX 6800 XT at ~$300 used, with its own ROCm headaches.

You think configuring SYCL backends and debugging driver issues sounds like a fun weekend. The A770 rewards that effort with usable AI performance at a bargain price.

You already run Linux with an Intel CPU and want to keep your stack consistent.

## Don't buy the A770 if…

You want NVIDIA-level simplicity. Install Ollama, run models, done. That's the NVIDIA experience. Intel requires more setup, more troubleshooting, more patience.

You're on Windows. While IPEX-LLM works on Windows, the ecosystem is strongest on Linux. Windows users should buy NVIDIA.

Image generation is your thing. The A770 is roughly half the speed of an RTX 3060 for Stable Diffusion. The extra VRAM doesn't compensate for the speed loss.

You need to fine-tune models. Training is effectively CUDA-only for consumer cards. The A770 is inference-only in practice.

Your time is worth more than $70. That's roughly the savings over an RTX 3060 12GB, and it comes with real software friction. If an afternoon of troubleshooting costs you more than that, buy the 3060.

## The Verdict

The Arc A770 16GB occupies a narrow but real niche: maximum VRAM at minimum price. At ~$255 used, nothing else gives you 16GB. That lets you run models that 8GB and 12GB cards simply cannot load.

The software is usable but unpolished. IPEX-LLM's Ollama integration is the easiest path and works today. llama.cpp SYCL is the safest long-term bet since it's merged into the official codebase. Both are slower than NVIDIA at the same model size, but fast enough for interactive chat.

My recommendation: **buy the A770 16GB used at ~$255 if you specifically need 16GB of VRAM and can't afford the $300+ alternatives.** The 4GB VRAM advantage over an RTX 3060 12GB lets you run 14B models that the 3060 can't fit. That capability gap is worth the software friction, if you're the kind of person who reads articles like this one.

If you just want things to work, buy an RTX 3060 12GB for $180 used and enjoy the CUDA ecosystem. No shame in that.

> Not sure what fits in 16GB of VRAM? Check the VRAM Planning Tool.

## Related Guides

- GPU Buying Guide for Local AI
- AMD vs NVIDIA for Local AI: Is ROCm Ready?
- What Can You Run on 16GB VRAM?
- Best Used GPUs for Local AI in 2026

Source: https://insiderllm.com/guides/intel-arc-local-ai/

Free guides for running AI locally