

## NVIDIA H100 PCIe Inference Power Profile

Public data on H100 PCIe power use are limited, but several sources offer insight. Measurements of real LLM inference on H100 show **per-query energy on the order of 0.1–0.3 Wh**. For example, a study found that generating a “thank you” response on LLaMA-3 (8B) used about 0.202 Wh on the GPU <sup>1</sup>. Likewise, an OPT-13B model consumed only **0.08 Wh** per query on one H100 (vs. 0.18 Wh on two L4 GPUs) <sup>2</sup>. A recent analysis projects a median of **0.34 Wh per query** on an H100 node for very large models <sup>3</sup>. These translate to *hundreds of watts* of GPU draw during inference: for example, ~0.34 Wh in a few seconds corresponds to roughly **250–300 W** average draw, i.e. ~70–80% of the card’s TDP. In short, the **steady-state inference draw** on H100 PCIe is validated to be on the order of a few hundred watts, consistent with our prior 250–280 W estimate ( $\approx$ 70–80% of 350/400 W TDP).

We found no direct public measurements for the **idle or phase-specific draws**. However, analogous data suggest an idle draw of only a few dozen watts. For instance, idle A100 GPUs draw ~50 W <sup>4</sup>, so the H100 PCIe likely idles on the order of **30–60 W**. This supports retaining our assumption of ~60–80 W idle for H100 PCIe. The card’s **peak power (TDP)** is 400 W for the H100 NVL (PCIe) variant <sup>5</sup> (standard PCIe card is 350 W). Actual peaks will not exceed these limits under heavy workload. We found no published **model-loading, warmup, or launch power** traces; in the absence of data we keep our prior estimates (e.g. model loading  $\approx$ 60–70% of inference power, warmup near inference power, launch ~30–50%) as plausible placeholders.

Studies of multi-GPU clusters confirm that **GPU power dominates system power** and tends to move in sync. One analysis notes a “strong correlation between GPU peak power and cluster-wide peak consumption,” highlighting that GPUs are the main energy contributors <sup>6</sup>. This qualitative finding is consistent with assuming a high correlation coefficient (C) across GPUs. In practice, our assumed C (0.3–0.7) is reasonable, though real synchronous workloads may show even higher correlation.

**Summary:** Extracted data validate our steady-state inference power (~250–280 W per H100) and peak/TDP values (350–400 W) <sup>1</sup> <sup>3</sup> <sup>5</sup>. Idle draw is likely only a few tens of watts (comparable to A100’s ~50 W) <sup>4</sup>. Phase transition powers remain unmeasured but are assumed somewhat lower (launch) or similar (load/warmup) to inference power. GPUs clearly dominate cluster power <sup>6</sup>, so our multi-GPU scaling assumptions appear sound. In short, H100 PCIe inference draws hundreds of watts during active inference, with idle draw tens of watts, and these findings align with published measurements <sup>1</sup> <sup>3</sup> <sup>2</sup> <sup>6</sup> <sup>5</sup> <sup>4</sup>.

**Citations:** Empirical measurements of H100 inference energy (converted above to power) come from Hugging Face benchmarks and recent analyses <sup>1</sup> <sup>3</sup> <sup>2</sup>. Idle reference is from NVIDIA forum data on A100 <sup>4</sup>. The H100 PCIe TDP (400 W) is given in NVIDIA’s product brief <sup>5</sup>. GPU-vs-cluster power correlation is noted in the literature <sup>6</sup>. All cited sources are open/public.

---

<sup>1</sup> Saying Thank You to a LLM Isn’t Free — Measuring the Energy Cost of Politeness  
<https://huggingface.co/blog/jdelavande/thank-you-energy>

- 2 6 From Prompts to Power: Measuring the Energy Footprint of LLM Inference  
<https://arxiv.org/html/2511.05597>
- 3 [2509.20241] Energy Use of AI Inference: Efficiency Pathways and Test-Time Compute  
<https://arxiv.org/abs/2509.20241>
- 4 Minimizing A100 power consumptin while idle under Linux - System Management and Monitoring (NVML) - NVIDIA Developer Forums  
<https://forums.developer.nvidia.com/t/minimizing-a100-power-consumptin-while-idle-under-linux/308388>
- 5 nvidia.com  
[https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/h100/PB-11773-001\\_v01.pdf](https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/h100/PB-11773-001_v01.pdf)