

# 188GB HBM3 Nvidia H100 NVL vznikla na míru pro Chat-GPT

[diit.cz/clanek/188gb-hbm3-nvidia-h100-nvl-vznikla-na-miru-pro-chat-gpt](https://diit.cz/clanek/188gb-hbm3-nvidia-h100-nvl-vznikla-na-miru-pro-chat-gpt)



Zdroj: Nvidia

Nvidia se pochlubila novou verzí akcelérátoru Hopper, která vznikla na míru pro Chat-GPT. Zajímavá je teoreticky až 2× vyšší energetickou efektivitou a nezvyklou kapacitou paměti, 2× 94 GB HBM3...

Nvidia loni uvedla akcelérátory Hopper ve dvou základních podobách. Jednak jako SXM5 moduly a jednak jako PCIe karty. Jako obvykle dosahuje SXM5 varianta vyššího výkonu i TDP, neboť je na modulu podstatně více prostoru pro velký radiátor chladiče. Nová varianta označená jako NVL má však opět podobu PCIe karty, respektive karet. Nvidia H100 NVL jsou totiž dvě PCIe karty spojené můstkem NVLINK. Nejde však o původní H100 PCIe. Nvidia k tomuto řešení sáhla patrně pro vyšší výpočetní denzitu takového řešení. Zkrátka méně prostoru zabraného chladičem ~ prostor pro více GPU. Jak vyplývá z parametrů, energetická efektivita má být 2× vyšší (2× více výkonu při stejném TDP díky dvojici karet s TDP polovičním oproti SXM5 provedení). Lze však očekávat, že realita bude trochu níže, protože reálné a papírové takty se budou patrně nezanedbatelně lišit.

	<b>Nvidia A100</b>	<b>Nvidia H100</b>		
<b>GPU</b>	GA100	GH100		2× GH100 „NVL“
<b>architektura</b>	Ampere	Hopper		
<b>formát</b>	SXM4	SXM5	PCIe	2× PCIe
<b>CU/SM</b>	108	132	114	2× 132?
<b>FP32 jader</b>	6912	15872 16896	14592	2× 16896?
<b>FP64 jader</b>	3456	8448	7296	2× 8448?
<b>INT32 jader</b>	6912	8448	7296	2× 8448?
<b>Tensor Cores</b>	432	528	456	2× 528?
<b>takt</b>	1410 MHz	1980 MHz	1750 MHz	1980 MHz?
↓↓↓ T(FL)OPS ↓↓↓				
<b>FP16</b>	78	120 134	102	2× 134
<b>BF16</b>	39	120 134	102	2× 134
<b>FP32</b>	19,5	60 67	51	2× 67
<b>FP64</b>	9,7	30 34	26	2× 34
<b>INT4</b>	?	?	?	?
<b>INT8</b>	?	?	?	?
<b>INT16</b>	?	?	?	?
<b>INT32</b>	19,5	30 34	26	2× 34
<b>FP8 tensor</b>	<b>X</b>	<u>1979/3958*</u>	<u>1513/3026*</u>	2× <u>1979/3958*</u>
<b>FP16 tensor</b>	312/624*	<u>989/1979*</u>	<u>757/1513*</u>	2× <u>989/1979*</u>
<b>BF16 tensor</b>	312/624*	<u>989/1979*</u>	<u>757/1513*</u>	2× <u>989/1979*</u>
<b>FP32 tensor</b>	19,5	60? 67?	51?	2× 67?
<b>TF32 tensor</b>	156/312*	495/989*	378/757*	2× 495/989*
<b>FP64 tensor</b>	19,5	67	51	2× 67
<b>INT8 tensor</b>	624/1248*	1979/3958*	1513/3026*	2× 1979/3958*
<b>INT4 tensor</b>	1248/2496*	?	?	?

	↑↑↑ T(FL)OPS ↑↑↑				
<b>TMU</b>	432		528	456	2× 528
<b>LLC</b>	40 MB		50 MB		2× 50 MB?
<b>sběrnice</b>	5120bit		5120bit		6144bit
<b>paměť</b>	40 GB	80 GB	80 GB		2× 94 GB
<b>HBM</b>	2,43 GHz	3,2 GHz	HBM3 5,23 GHz	HBM2E 3,2 GHz	HBM3 5,1 GHz
<b>pam. prop.</b>	1555 GB/s	2048 GB/s	3350 GB/s	2048 GB/s	2× 3,9 TB/s
<b>TDP</b>	400 W		700-800 W	350 W	700 W
<b>transistorů</b>	54,2 mld.		80 mld.		2× 80 mld.
<b>plocha GPU</b>	826 mm <sup>2</sup>		814 mm <sup>2</sup>		2× 814 mm <sup>2</sup>
<b>proces</b>	7 nm			4nm	
<b>datum</b>	5. 2020	11. 2020	2022?		H2 2023

Protože jsou jazykové modely Chat-GPT velké, osadila Nvidia poprvé GPU GH100 všemi šesti HBM moduly (tzn. plně využila 6144 bit sběrnici) a použila navíc HBM3 (oproti HBM2E na původní PCIe variantě). To by teoreticky odpovídalo  $6 \times 2 \times 16 \text{ GB} = 192 \text{ GB}$  paměti, jenže reálně specifikace uvádí 188 GB. Jak toho lze dosáhnout, když žádné 15,66GB moduly HBM neexistují? Nvidia se nejspíš dohodla s některým výrobcem pamětí, který jí dodává moduly s deaktivovanými vadnými buňkami, což by mohlo znamenat mírně nižší kapacitu, ale i podstatně výhodnější cenu.

Přestože Nvidia hovoří o vydání H100 NVL, reálně mají být tyto akcelerátory dostupné někdy ve druhém pololetí letošního roku.

[nahlásit chybu](#)



### **Nvidia nechala redaktory ukázat GeForce RTX 4090 [fotografie]**

---

GeForce RTX 4090 sice bude vydána až příští týden (zároveň se zahájením prodejů Intel Arc A7x0), ale Nvidia nechala redakce zveřejnit fotografie karty již nyní, v podstatě zároveň s recenzemi Arc...



## GeForce RTX 4080 v GeekBench: 10 % nad RTX 3090 Ti, 34 % nad RTX 3080

16GB (tedy ta nezrušená) verze GeForce RTX 4080 byla otestována v GeekBench 5. Její výkon dosahuje v průměru 10 % nad GeForce RTX 3090 Ti a 34 % nad GeForce RTX 3080...



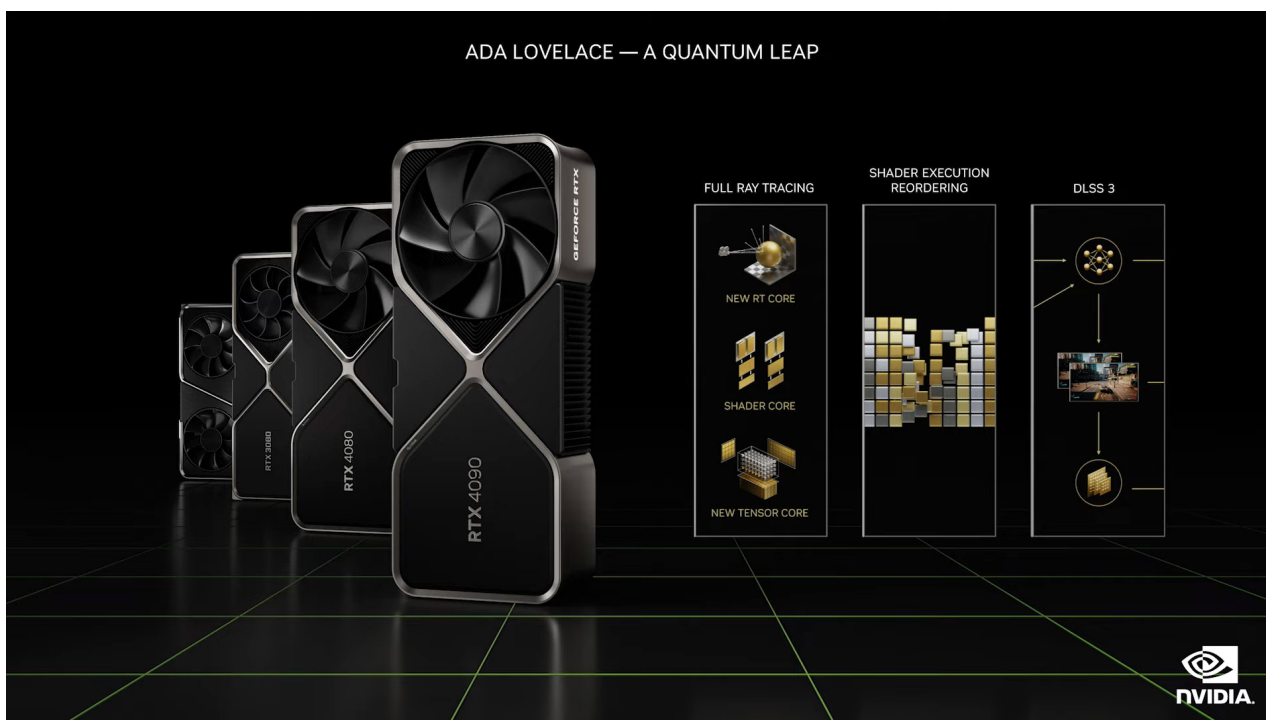
Uživatelé: GeForce RTX 4080 by měla stát nanejvýš \$800

Z uživatelské ankety, které se zúčastnilo téměř 12 tisíc respondentů, vyplynulo, že 97,8 % uživatelů považuje GeForce RTX 4080 za předraženou, pouze 0,9 % se s její cenou ztotožňuje...



## **GeForce RTX 4090 v GeekBench: 60-63 % nad GeForce RTX 3090 Ti**

Skóre GeForce RTX 4090, nové generace grafického high-endu od Nvidie, se objevilo v online databázi benchmarku GeekBench...



## GeForce RTX 4080 16GB se původně měla jmenovat GeForce RTX 4070

Nakonec vše nasvědčuje tomu, že původní plán Nvidie byl podobný jako u generace Ampere. Na největším jádru řady 102 postavit GeForce RTX x090 a o stupeň menší řadu 103 použít pro GeForce RTX x070...



**Jiří "no-X" Souček**

[více článků, blogů a informací o autorovi](#)

## Diskuse ke článku 188GB HBM3 Nvidia H100 NVL vznikla na míru pro Chat-GPT

před 1 dnem | Wladows | [Poběží na tom i Skynet?](#)

před 2 dny | no-X | [To jsou pěkné teorie, které by dávaly smysl, ale...](#)

před 2 dny | melkor | [A pak, že Sci-Fi je brak. 87 let od vydání a...](#)

před 2 dny | TyNyT | [a tohle sedí taky: Přes tento střízlivý úsudek...](#)

před 2 dny | TyNyT | [Na konci protokolu shrnula odborná komise...](#)

před 2 dny | jk2 | Dle mých zkušeností je nutné výsledky z Chat-GPT...

před 2 dny | danieel | Ta pamet bude mit spise chybejici kanal. Kazdy...

před 2 dny | RedMaX | Ptali se Chat-GPT a říkal, že je současný český...

Zobrazit diskusi

Pro psaní komentářů se, prosím, přihlaste nebo registrujte.