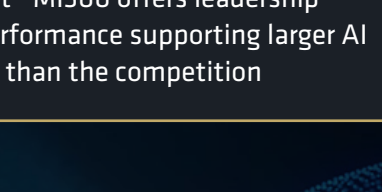




PCS

9

<b>AMD EPYC</b>	• Small to Medium Models • Medium to Large Models	<b>AMD INSTINCT</b>	• Small to Medium Models • Batch/Small Scale Inference
---------------------	--	-------------------------	---



## ARM® CORTEX™ PROCESSORS

consolidation  
ready for AI, a  
the AI inference

PERFORMANCE UPLIFT AND ENERGY

AMD Ryzen™ AI is the leading CPU for AI P with a wide variety of laptop, desktop, and

Microsoft  
+ Procyon Office Pro



95%

Proven Off

Less Power

6242

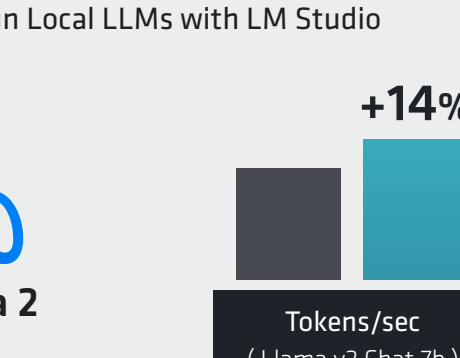
Tel: +33 239 41 00 00

D

## STUDY IN THE WORKED

## THE PC FLEET

See Endnotes: HWKP-26, HWKP-27, HWKP-28, HWKP-29



Model	Time to 1st Token
Baseline	~1.5s
Ours	~2.5s

Time to 1st Token  
(llama-2-chat-7b)

CPU (20W) AMD Ryzen™ 3 PRO 3200U CPU (15W)

See Endnote: PUY-58

CloudEra Couchbase DataStax elastic MarkLogic mongoDB splunk>

[Home](#)
[About Us](#)
[Contact Us](#)
[Privacy Policy](#)

# Extending From “Day 0” to “Bleeding-Edge” Support



Incredible performance upgrade for  
GenAI training and inference

Technologies  
Microsoft

© 2006 The Authors  
Journal compilation © 2006 Blackwell Publishing Ltd

---

Outlook, Teams + Procyon Office Productivity Power Point, Teams + Procyon Office Productivity Word

The Intel based server power cost of \$652.844, using a PUE of 1.7, saving \$168.719 over the 3 years of this analysis with an estimated US power cost of \$0.128 / kWh. The ZP EPC core CPU solution also provides estimated Greenhouse Gas Emission savings emissions avoided equivalent to 1.075 MTCO<sub>2</sub> (1185 US tons) over the 3 years of this analysis which is 395 US tons of CO<sub>2</sub> annual and is the equivalent of the sequestration equivalent of 2430 US tons annually.

M330-39: Number of simultaneous text generating copies of the Llama2-70b chat model, using vLLM, comparison using custom docker container for each system based in AMD internal testing as of 1/16/2023. Configurations: 2P Intel Xeon Platinum 8480C CPU server with 8x AMD Ryzen™ (9526h, 7590W) CPUs, ROCm™ v5.0 pre-release, PyTorch 2.0.2, vLLM for ROCm, Ubuntu 22.04.2. Vs. on Nvidia DGX H100 with Intel Xeon Platinum 8480C Processors, 8x Nvidia H100 (80GB, 700W) CPUs, CUDA 12.1, PyTorch 2.1.0, vLLM v0.2.2 (most recent), Ubuntu 22.04.3. Performance may vary based on use of latest drivers and optimizations.

M330-40: Testing completed 1/16/2023 by AMD performance lab using MosaicML, yielding difference to the fine tune the MPT-30B model for 2 epochs using the MosaicML instruct-v1 dataset and a max sequence length of 8192 tokens using custom docker container for each system

to Nvidia HX100 (B0C08, 700W) GPU, CUDA 11.8, PyTorch 2.0.1, MosaicML, lin-foundation, Ubuntu 22.04.3. Server manufacturers may use different configurations, yielding different results. Performance may vary from latest drivers and optimizations.

PHX-59. Testing as of Feb 2025 by AMD. Sustained performance average of multiple runs with specimen prompt "Write me a story about an orange cat called mr whiskers". All tests conducted on LM Studio 0.2.16. Performance may vary. Market price retrieved on 3/4/2025 (Amazon, US).  
 Phoenix HP Pavilion Plus Laptop 14 w/ Ryzen 7 7840U 15W, 16GB LPDDR5 6400, Windows 23H2 22631.355, Adrenalin Driver 24.2.1. Meteor Lake: Acer Swift SFC14-72T, Intel Core Ultra i7 155H 28W 16GB LPDDR5 6400, Windows 23H2 22631.355, Distro 21.0.101.5333.

Ryzen AI is defined as the combination of a dedicated AI engine, AMD Radeon™ graphics engine, and Ryzen processor cores that enable AI capabilities. OEM and ISV enhancement is required, and certain AI features may not yet be optimized for Ryzen AI processors. Ryzen AI

compatible with: (a) AMD Ryzen 7040 and 8040 Series processors except Ryzen 5 7540U, Ryzen 5 8540U, Ryzen 3 7440U, and Ryzen 3 8440U processors; and (b) all AMD Ryzen 8000C Series desktop processors except the Ryzen 5 8500G/GE and Ryzen 3 8300G/CE. Please check with your system manufacturer for feature availability prior to purchase. CG-020b.

© 2024 Advanced Micro Devices, Inc. all rights reserved. AMD, the AMD arrow, AMD Instinct, EPYC, Radeon, RDCM, Ryzen, Threadripper, and combinations thereof, are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective owners.