

## USER GUIDE AMD EPYC 9004



Publication Revision Issue Date 58736 1.1 August, 2024

#### © 2024 Advanced Micro Devices, Inc. All rights reserved.

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information. Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale. GD-18u

#### Trademarks

©2024 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD logo, EPYC, and combinations thereof are trademarks of Advanced Micro Devices. Nutanix is a trademark or registered trademark of Nutanix in the US or other countries. Certain AMD technologies may require third-party enablement or activation. Supported features may vary by operating system. Please confirm with the system manufacturer for specific features. No technology or product can be completely secure.

\* Links to third party sites are provided for convenience and unless explicitly stated, AMD is not responsible for the contents of such linked sites and no endorsement is implied.

Date	Version	Changes
Aug, 2024	1.0	Initial public release
Aug, 2024	1.1	Added model description and use case information; minor errata fixes

# **Table of Contents**

Introduction & Requirements		
About Large Language Models Use Cases		
1.2.1 Boost Employee Experience	2	
1.2.2 Transform Customer Interactions	2	
1.2.3 Enhance Internal Collaboration	3	
Technologies	3	
AMD EPYC 9004 Series Processors		
Important Reading		
Requirements		
Caveats	4	
Host BIOS Settings	5	
Performance Optimization	9	
.1 AVX Configuration		
EPYC NUMA Alignment	10	
Deploying Quantized Llama2 on Nutanix Cluster (TorchServe kind)	13	
Build a Llama2 TorchServe Model Archive		
4.1.1 Step 1: Setup PVC and Create a Temporary Pod		
4.1.2 Step 2: Install the tools and model	14	
4.1.3 Step 3: Modify Model Configuration File and Create the Model Archive	14	
4.1.4 Step 4: Copy Required Files to the model-store-pod	15	
4.1.5 Step 5: Copy model_store Contents	15	
Deployment	15	
Benchmarking Llama2 Inference Serving on AMD EPYC Systems	19	
Explaining Benchmarking	19	
Locust		
5.2.1 Benchmark Client	21	
Linear Llama2 7B "CPU-Only" Scaling	23	
	Introduction & Requirements         About Large Language Models         Use Cases         1.2.1       Boost Employee Experience         1.2.2       Transform Customer Interactions         1.2.3       Enhance Internal Collaboration         Technologies       AMD EPYC 9004 Series Processors         Important Reading       Requirements         Caveats       Requirements         AVX Configuration       AVX Configuration         EPYC NUMA Alignment       Performance Optimization         Build a Llama2 TorchServe Model Archive       4.1.1         Step 1: Setup PVC and Create a Temporary Pod       4.1.2         4.1.2       Step 2: Install the tools and model         4.1.3       Step 3: Modify Model Configuration File and Create the Model Archive         4.1.4       Step 4: Copy Required Files to the model-store-pod         4.1.5       Step 5: Copy model_store Contents         Deployment       Benchmarking Llama2 Inference Serving on AMD EPYC Systems	

## 

Chapter

# Introduction & Requirements

This document guides you through the procedures involved in configuring and setting up a Nutanix cluster with a focus on deploying LLama2 models and performance benchmarking of multi-user requests that specifically evaluate Large Language Models (LLMs), such as Ilama-2-7b-int8 and its variant models. It contains detailed, step-by-step instructions for establishing an Nutanix Kubernetes Engine (NKE) multi-node cluster environment, deploying LLama2, and conducting performance benchmarking. It also addresses the deployment process on both Nutanix VMs and cluster and specifies the minimum hardware and software configurations necessary for this setup.

LLM-in-a-Box is a turnkey AI solution for organizations wanting to implement GPT capabilities while maintaining control of their data and applications. It includes everything needed to build AI-ready infrastructure, including:

- Nutanix<sup>™</sup> Cloud Platform infrastructure.
- Nutanix Files and Object storage for running and fine-tuning GPT models.
- Open-source software to deploy and run AI workloads, including PyTorch, and Kubeflow.
- Support for a curated set of LLMs, including Llama2, Falcon, and MPT.

## 1.1 About Large Language Models

Large Language Models (LLMs) are revolutionizing the world and reshaping how we interact with technology. Models powered by advanced artificial intelligence, such as GPT-3 and Llama2, enable breakthroughs in natural language processing by comprehending and generating human-like text. LLMs are accelerating innovation in diverse applications from crafting creative content to enhancing customer support. Their ability to understand context, generate coherent responses, and adapt to various tasks has opened new frontiers in automation, education, and research. LLMs are driving transformative change by promising a future where communication and problem-solving reach unprecedented levels of efficiency and sophistication.

LLMs offer incredible potential for automating tasks, enhancing customer experiences, and driving research, but their reliance on expensive GPUs often creates a barrier to entry. This User Guide presents a groundbreaking solution that makes cutting-edge AI accessible to any organization deploying and serving Llama2 on systems powered by 4th Gen AMD EPYC<sup>™</sup> processors, allowing you to unlock the benefits of LLMs on your existing infrastructure. This solution offers the following benefits:

- Eliminates the need for expensive GPUs and enhances access to AI technology.
- Enhances Llama2 performance for fast text generation, natural language processing, and content creation by using 4th Gen AMD EPYC processors.
- Streamlines Llama2 deployment on the popular Red Hat<sup>®</sup> OpenShift<sup>®</sup> Container Platform (OCP) for seamless scalability, flexibility, and efficient resource management.
- Measure performance using detailed performance benchmarks.

Some of the ways you can use this solution include:

- Automating tasks to streamline workflows, boost efficiency, and empower your workforce by automating repetitive tasks, such as data analysis and report generation.
- Enhancing customer service by delivering personalized experiences through intelligent chatbots, sentiment analysis, and AI-powered recommendations.
- Driving research efforts with LLM-powered automated data analysis, document summarization, and hypothesis generation.
- Enhancing competitiveness by leveraging AI-driven opportunities for innovation across text and language-driven domains.

## 1.2 Use Cases

This section briefly explores some compelling AI use cases specifically designed for internal and external applications within your company. Artificial intelligence (AI) is rapidly transforming business landscapes, offering exciting opportunities to enhance efficiency, improve customer experience, and drive innovation. The synergy between AMD EPYC processors and OpenShift AI presents a compelling opportunity for businesses to unlock the transformative power of artificial intelligence. The combination of exceptional performance, optimized scalability, and robust platform support empowers you to tackle diverse AI workloads efficiently and cost-effectively. Combining AMD EPYC processors and OpenShift AI provides a solid foundation for your AI journey and keeps you ahead of the curve as technology continues evolving.

#### 1.2.1 Boost Employee Experience

- **Streamlined onboarding:** Imagine an AI assistant driving a smooth and efficient onboarding experience by guiding new hires through company policies, answering questions, and providing personalized knowledge access.
- **Instant support:** Picture employees saving time and reducing reliance on support teams by easily resolving common IT issues through intuitive AI-powered troubleshooting.
- **Seamless knowledge transfer:** Envision a system that fosters collaboration and self-service for immediate answers by seamlessly connecting employees to relevant internal knowledge bases.

### 1.2.2 Transform Customer Interactions

- **24/7 customer care:** Imagine an AI system ensuring consistent support around the clock by acting as your first line of defense by handling basic inquiries and freeing up human agents for complex issues.
- Intelligent lead nurturing: Picture an AI assistant freeing up sales reps for high-value interactions and closing deals by pre-qualifying leads, scheduling appointments, and nurturing potential customers.
- **Empowering self-service:** Envision customers readily accessing product information, resolving simple issues, and finding answers independently via an AI-powered support system.



#### 1.2.3 Enhance Internal Collaboration

- Effortless meeting management: Imagine an AI assistant boosting team alignment and follow-through by automatically summarizing meetings, generating action items, and distributing those action items.
- Enhanced project management: Picture an AI system facilitating seamless communication and collaboration across teams by keeping track of tasks, deadlines, and resources.
- Sparking innovation: Envision an AI assistant unlocking your team's creative potential by fueling brainstorming sessions with insightful suggestions, capturing ideas, and documenting outcomes.

### 1.3 Technologies

4th Gen AMD EPYC<sup>™</sup> processors are ideal for AI workloads because of their high core counts and strong performance. Nutanix Kubernetes Engine (NKE) simplifies deploying and managing machine learning pipelines, and Kserve optimizes model serving for fast predictions. The Llama2 LLM brings advanced natural language processing capabilities, and Nutanix AOS allows centralized management and monitoring of the entire AI infrastructure. This combination unlocks transformative AI experiences on your existing infrastructure.

## 1.4 AMD EPYC 9004 Series Processors

AMD EPYC 9004 Series Processors continue to redefine the standards for modern datacenters. 4th Gen AMD EPYC processors are built on the innovative x86 architecture and "Zen 4" core. 4th Gen AMD EPYC processors deliver efficient, optimized performance by combining high frequencies, the largest-available L3 cache, up to 128 (1P) or up to 160 (2P) lanes of PCle<sup>®</sup> Gen 5 I/O, synchronized fabric and memory clock speeds, and support for up to 6 TB of DDR5-4800 memory. Built-in security features, such as AMD Infinity Fabric<sup>™</sup> technology, Secure Memory Encryption (SME), and Secure Encrypted Virtualization (SEVSNP) help protect data while it is in use. AMD Infinity Guard features vary by EPYC<sup>™</sup> Processor generations and/or series. (Infinity Guard security features must be enabled by server OEMs and/or Cloud Service Providers to operate. Check with your OEM or provider to confirm support of these features. Learn more about Infinity Guard at http://www.amd.com/en/products/processors/server/epyc/infinity-guard.html. GD-183A.)

## 1.5 Important Reading

Please be sure to read the following guides (available from the <u>AMD Documentation Hub</u>), which contain important foundational information about 5th Gen AMD EPYC processors:

- AMD EPYC<sup>™</sup> 9005 Processor Architecture Overview
- BIOS & Workload Tuning Guide for AMD EPYC™ 9005 Series Processors

## 1.6 Requirements

The following requirements must be met in order to deploy LLM-in-a-Box solutions on Nutanix AOS/AHV:

- Nutanix cluster with AOS/AHV.
- A Prism Central installed on the same cluster (required for Nutanix Kubernetes Engine [NKE]).
- A Kubernetes cluster initialized using NKE, where each worker node is configured with 76 vCPUs, 128 GB memory, and 1 vNUMA node.

### 1.7 Caveats

Obtaining maximum performance requires fine tuning the software packages described in this user guide. Some of the procedures described herein may not be standard for a typical production environment. This specifically applies to the following operations that use the foundational Linux supports instead of Nutanix AOS management facilities:

- The Kuberenets worker VMs must be manually configured to use host CPU passthrough.
- The Kubernetes worker VMs must be manually pinned to the desired CPU cores.
- The Kubernetes worker VMs must be manually pinned to the desired NUMA node.



Chapter

# **Host BIOS Settings**

Table 2-1 lists the recommended host BIOS settings for LLM-in-a-Box on a Nutanix cluster. You must power cycle the host after modifying BIOS settings.

Name	Recommended Setting	Description	
Global C-State Control	Auto	<ul> <li>Enabled/Auto: Controls IO based C-state generation and DF C- states, including core processor C-States</li> </ul>	
		<ul> <li>Disabled: AMD strongly recommends not disabling this option because this also disables core processor C-States.</li> </ul>	
DF C-States	Disabled	Controls DF C-states.	
		<ul> <li>Disabled: Prevents the AMD Infinity Fabric from entering a low- power state.</li> </ul>	
		<ul> <li>Enabled/Auto: Allows the AMD Infinity Fabric to enter a low- power state.</li> </ul>	
Power Profile Selection	Auto	Auto/0: High-performance mode	
		• 1: Efficiency mode	
		2: Maximum I/O performance mode	
Core Performance Boost	Auto	Enabled/Auto: Enables Core Performance Boost.	
		Disabled: Disables Core Performance Boost.	
Determinism Control	Manual	Auto: Use default performance determinism settings.	
		Manual: Specify custom performance determinism settings.	
Determinism Enable	Auto	Auto/0: Power.	
		• <b>1:</b> Performance.	

Table 2-1: Recommended host BIOS settings for deploying LLM-in-a-Box on a Nutanix cluster

Nodes Per Socket (NPS)	1	<ul> <li>Memory Interleaving: The NPS setting always determines the memory interleaving regardless of whether LLC as NUMA is Enabled or Disabled.</li> <li># of NUMA nodes (if LLC as NUMA Domain is Disabled):</li> <li>NPS1: One NUMA node per socket (Most cloud providers use this as it provides consistent average memory latency to all the accesses within a socket).</li> <li>NPS2: Two NUMA nodes per socket.</li> <li>NPS4: Four NUMA nodes per socket</li> <li>NPS0 (not recommended): Only applicable for dual-socket systems. A single NUMA node is created for the whole two-socket platform.</li> <li>AMD recommends either NPS1 or NPS4 depending on your use case.</li> </ul>
		<b>Windows systems:</b> Make sure that the number of logical processors per NUMA node is <=64. You can do this by using NPS2 or NPS4 instead of the default NPS1.
TSME	Auto	Auto/Disabled: Disables transparent secure memory encryption.
		Enabled: Enables transparent secure memory encryption.
SEV Control	Disabled	In a multi-tenant environment (such as a cloud), Secure Encrypted Virtualization (SEV) mode isolates virtual machines from each other and from the hypervisor.
		Disabled: SEV is disabled.
		• Enabled: SEV IS enabled.
		If you disable and then reenable SEV, then you will need to power cycle your system after changing this setting back to <b>Enabled</b> .
SEV-ES	Disabled	Secure Encrypted Virtualization-Encrypted State (SEV-ES) mode extends SEV protection to the contents of the CPU registers by encrypting them when a virtual machine stops running. Combining SEV and SEV-ES can reduce the attack surface of a VM by helping protect the confidentiality of data in memory.
		Disabled: SEV-ES is disabled.
		Enabled: SEV-ES is enabled.

Table 2-1: Recommended host BIOS settings for deploying LLM-in-a-Box on a Nutanix cluster (Continued)

SEV-SNP Support	Disabled	Secure Encrypted Virtualization-Secure Nested Paging (SEV-SNP) mode builds on SEV and SEV-ES by adding strong memory integrity protection to create an isolated execution environment that helps prevent malicious hypervisor-based attacks such as data replay and memory re-mapping. SEV-SNP also introduces several additional optional security enhancements that support additional VM use models, offer stronger protection around interrupt behavior, and increase protection against recently-disclosed side channel attacks.
		• <b>Disabled:</b> SEV-SNP is disabled.
		• Enabled: SEV-SNP is enabled.
1		

Table 2-1: Recommended host BIOS settings for deploying LLM-in-a-Box on a Nutanix cluster (Continued)

Chapter

# **Performance Optimization**

The GPT stack is deployed in a pod running within a K8s node. NKE runs each K8s node in a VM. Obtaining maximum performance requires configuring the K8s VM node to utilize the AMD EPYC processor AVX features and to align with AMD EPYC processor NUMA topology. This chapter refers to both single and multiple K8s node VM(s) as VMk8s-node.

## 3.1 AVX Configuration

As mentioned above, all VMk8s-node must be configured to use the AMD EPYC AVX2 and AVX512 instruction sets. By default, the Nutanix AHV does not support 4th Gen AMD EPYC processor, but you can work around this by configuring the hypervisor to passthrough the host CPU type to VMk8s-node via the cpu passthrough AHV VM attribute.

You should also make sure that VMk8s-node is not running two AVX threads on the same CPU core simultaneously because this will cause AVX resource contention on the CPU. You can avoid this by configuring VMk8s-node to run with a single thread using the num\_threads\_per\_core AHV VM attribute.

To enable cpu\_passthrough and num\_threads\_per\_core for a VMk8s-node:

- 1. Login to any of the CVMs in the cluster as the nutanix user.
- 2. Obtain the list of all VMs running on the cluster. acli vm.list
- 3. Power off the VMk8s-node. acli vm.off <vm name>
- 4. Set the VM attributes, acli vm.update <vm\_name> cpu\_passthrough="true" num\_threads\_per\_core="1"
- 5. Power on the VMk8s-node acli vm.on <vm name>
- 6. Repeat Steps 3-5 for all VMk8s-node running the GPT stack.

## 3.2 EPYC NUMA Alignment

The VMk8s-node must be configured to run within an AMD EPYC processor NUMA node. To do this:

- 1. On Prism Central, migrate the NKE related VMs that are not VMk8s-node to one of the AHV hosts in the cluster.
- 2. Login to one of the CVMs, then power off all the VMk8s-node that run the GPT stack.
  - a. Obtain the list of all VMs running on the cluster. acli vm.list
  - b. Power off the VMk8s-node. acli vm.off <vm name>
- 3. Execute the following commands to set VMk8s-node to physical NUMA node affinity and to set the VM to host affinity.

```
acli vm.update <vm name> extra_flags="numa_pinning=1"
acli vm.affinity set <vm name> host list=<host id>
```

Note: Each VMk8s-node should have a distinct host affinity. Avoid sharing an AHV host with more than one VMk8snode.

 Execute the following command to power up each VMk8s-node. acli vm.on <vm name>

Follow Steps a-b on each AHV host that runs a VMk8s-node.

a. Execute the following command to ensure the VMk8s-node is running on NUMA node #1. If not, then power off the VM and then power it on again; repeat this as necessary until the VM is running on NUMA node #1.

#### Note: The numa pinning setting does not set the NUMA affinity 100% of the time.

numastat -p qemu-kvm

The output appears as shown below. Execute the ps command to find out which PID corresponds to VMk8s-node. The memory usage should appear in the Node 1 column.

Per-node process memory usage (in MBs)

	0		
PID	Node 0	Node 1	Total
14143 (qemu-kvm)	49197.79	4.78	49202.57
1467714 (qemu-kvm)	37.60	49165.80	49203.40
Total	49235.39	49170.58	98405.96

b. Execute the acli command on the CVM to find out the UUID of the VMk8s-node. ssh nutanix@<cvm-ip> /usr/local/nutanix/bin/acli vm.list | fgrep worker

#### For example:

# ssh nutanix@192.168.5.254 /usr/local/nutanix/bin/acli vm.list | fgrep worker-0 Nutanix Controller VM

karbon-nke-titanite-26e2b2-worker-0 4cc166c0-5f24-42c7-8c62-b0ae291506d2



c. Execute the following command to pin the vCPUs of VMk8s-node to a fix set of cores. virsh vcpupin <vm uuid> <vcpu #> <core #>

The following example assumes the VMk8s-node has 76 vCPUs, and the host CPU has 96 cores. The cores on the NUMA node #1 correspond to vCPU cores 96 to 191 on a 96-core 4th Gen AMD EPYC processor. It is good practice to pin the VMk8s-node starting at core number=<high core number>-<number vCPUs>+1. This example uses have 191-76+1=116.

virsh vcpupin 4cc166c0-5f24-42c7-8c62-b0ae291506d2 0 116 virsh vcpupin 4cc166c0-5f24-42c7-8c62-b0ae291506d2 1 117 virsh vcpupin 4cc166c0-5f24-42c7-8c62-b0ae291506d2 2 118 . . . virsh vcpupin 4cc166c0-5f24-42c7-8c62-b0ae291506d2 74 190 virsh vcpupin 4cc166c0-5f24-42c7-8c62-b0ae291506d2 75 191

### Chapter

# Deploying Quantized Llama2 on Nutanix Cluster (TorchServe)

This chapter describes how to deploy the Hugging Face Llama2-chat-hf int8 model on NKE.

### 4.1 Build a Llama2 TorchServe Model Archive

### 4.1.1 Step 1: Setup PVC and Create a Temporary Pod

1. Create a PVC by executing the following commands in the terminal:

```
$ kubectl apply -f pvc.yaml
## pvc.yaml ##
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
   name: model-store-claim
spec:
   resources:
    requests:
      storage: 80Gi
accessModes:
      - ReadWriteOnce
```

2. Execute the following commands to mount the PVC named model-store-pod on temp pod:

```
$ kubectl apply -f temp pod.yaml
## temp pod.yaml ##
apiVersion: v1
kind: Pod
metadata:
 name: model-store-pod
spec:
  volumes:
    - name: model-store
     persistentVolumeClaim:
       claimName: model-store-claim
  containers:
    - name: model-store
     image: ubuntu
     command: [ "sleep" ]
      args: [ "infinity" ]
      volumeMounts:
        - mountPath: "/mnt"
         name: model-store
      resources:
        limits:
```

## AMD, Nutanix<sup>™</sup> AOS/AHV AI-in-a-Box Solutions Build Guide

cpu: 2000m memory: 5Gi

#### 4.1.2 Step 2: Install the tools and model

- 1. Install dependencies and make a symbolic link inside the temp pod. \$ apt-get update && apt-get install git python3 pip \$ ln -s /usr/bin/python3 /usr/bin/python
- 2. Clone the <u>Torchserve</u>\* repository.
- 3. Execute the following commands:
  - \$ cd serve
  - \$ python3 ./ts\_scripts/install\_dependencies.py
  - \$ pip install torchserve torch-model-archiver \$ torch-workflow-archiver
- 4. Download the model.
  \$ cd examples/large\_models/Huggingface\_accelerate
  \$ apt-get install libopenmpi-deV
- 5. Modify requirements.txt by adding the following lines:
   \$ llama-cpp-python==0.2.78
   \$ pip install -r requirements.txt
   \$ pip install -U "huggingface\_hub[cli]"
   \$ huggingface-cli login
- 6. Execute the following command to download the desired model to the model folder: \$ python3 Download model.py --model path <path> --model name <huggingface-repo>

```
To download a quantized Llama2-7b-chat model:
$ python3 Download_model.py --model_path model --model_name arunimad/Llama2-7b-chat-
int8
```

### 4.1.3 Step 3: Modify Model Configuration File and Create the Model Archive

- 1. Modify model-config.yaml as follows:
  - Edit model\_name as desired. For example: llama-cpp-7b
  - Edit the model\_path variable to point to the quantized model. For example: model/models--arunimad--Llama2-7b-chat-int8/snapshots/ 4525405b7837e6d0a5fc98e295e9646942d33b4f/ggml-model-q8.gguf
- 2. Modify llama\_cpp\_handler.py by replacing the inference function with the following function:

```
def inference(self, data):
    prompt_template = '''[INST] <<SYS>>You are a helpful chatbot that will answer to the
prompt as needed.<</SYS>>{prompt}[/INST]'''
    result = self.model(prompt_template.format(prompt = data["prompt"]),
max_tokens=data["max_tokens"], top_p=data["top_p"],temperature=data["temperature"],
echo=False)
    tokens = self.model.tokenize(bytes(data["prompt"], "utf-8"))
    return result
```



3. Create the model archive using torch-model-archiver:

```
torch-model-archiver --model-name llamacpp-7b --version 1.0 --handler llama_cpp_handler.py
--config-file model-config.yaml --archive-format no-archive --requirements-file
requirements.txt
```

- 4. Create a model store folder.
- 5. Copy the archive folder into model\_store.
- 6. Copy model/folder with the downloaded model into the archive folder.

#### 4.1.4 Step 4: Copy Required Files to the model-store-pod

- 1. Execute into the pod.
- 2. Create model-store and config folders by executing the following commands:
  - \$ kubectl exec -it model-store-pod bash
  - \$ cd mnt
  - \$ mkdir model-store
  - \$ mkdir config
- 3. Copy the model files into the PVC pod.
   \$ kc cp model\_store/llama-cpp-7b model-store-pod:/mnt/model-store/
- 4. Add config.properties to the config folder.

```
$ kc cp config.properties model-store-pod:/mnt/config/
```

```
inference_address=http://0.0.0.0:8080
management_address=http://0.0.0.0:8081
metrics_address=http://0.0.0.0:8082
enable_envvars_config=true
install_py_dep_per_model=true
load_models=all
model_store=/home/model-server/model_store
model_snapshot={"name":"startup.cfg","modelCount":1,"models":{"llama-cpp-
7b":{"1.0":{"defaultVersion":true,"marName":"llama-cpp-
7b","minWorkers":1,"maxWorkers":1,"batchSize":1,"maxBatchDelay":200,"responseTimeout":1200
}}}
```

### 4.1.5 Step 5: Copy model\_store Contents

Copy the model\_store contents from serve/examples/LLM/llama/chat\_app/model\_store to /mnt/
model\_store.

## 4.2 Deployment

1. Create deploy\_int8.yaml.

```
apiVersion: apps/v1
kind: Deployment
metadata:
name: ts-def
labels:
app: ts-def
spec:
replicas: 1
selector:
```

## AMD Nutanix<sup>™</sup> AOS/AHV AI-in-a-Box Solutions Build Guide

```
matchLabels:
     app: ts-def
 template:
   metadata:
     labels:
       app: ts-def
   spec:
     volumes:
        - name: model-store
         persistentVolumeClaim:
           claimName: model-store-claim-int8
     containers:
        - name: torchserve
          image: pytorch/torchserve:latest-cpu
          command: ["/bin/sh", "-c"]
          args: ["torchserve --start --ts-config /home/model-server/config/
config.properties; sleep infinity"]
          env:
            - name: LOG LOCATION
             value: "/home/model-server/logs"
            - name: METRICS LOCATION
             value: "/home7model-server/logs"
            - name: TEMP
              value: "/home/model-server/tmp"
            - name: OPENBLAS NUM THREADS
             value: "1"
         ports:
            - containerPort: 8080
            - containerPort: 8081
            - containerPort: 8082
          volumeMounts:
            - name: model-store
             mountPath: /home/model-server
          resources:
            limits:
              cpu: 48
             memory: 48Gi
            requests:
             cpu: 48
             memory: 48Gi
          securityContext:
            allowPrivilegeEscalation: false
            runAsUser: 0
```

- 2. Execute the following command to deploy the TorchServe server: \$ kubectl apply -f deploy\_int8.yaml
- 3. Check the pod logs in the deployment. If there are any errors, then proceed to Step 4, else skip to Step 6. \$ kubectl logs pod-name
- 4. Create a tmp directory in the torchserve pod.
  - \$ kubectl exec -it pod-name bash
  - \$ cd /home/model-server
  - \$ mkdir tmp
  - \$ Exit pod and delete pod :
  - \$ kubectl delete pod pod-name
- 5. Check the new pod logs to verify that the model loads.



Execute the following commands to create a service: 6.

```
$ kubectl apply -f svc.yaml
apiVersion: v1
kind: Service
metadata:
 name: svc
 namespace: int8
spec:
  type: LoadBalancer
  selector:
   app: ts-def
 ports:
  - name: inference
   port: 8080
    targetPort: 8080
  - name: inference2
    port: 8085
    targetPort: 8085
  - name: mgmt
   port: 8081
   targetPort: 8081
  - name: metrics
   port: 8082
    targetPort: 8082
  - name: streamlit01
    port: 8501
    targetPort: 8501
```

You can now send requests using curl commands. Use the following templates, as appropriate:

#### LLM Endpoint:

http://<routename>/predictions/<model-name>

#### Sample Endpoint: ٠

```
curl -v "http://localhost:8080/predictions/llamacpp-7b" -H 'Content-Type: application/
json' -d "@sample.json"
```

#### Sample JSON: ٠

```
"prompt": "1+1",
"max_tokens": 128,
"top p": 1.0,
"temperature": 1.0
```

## 

## Chapter

# Benchmarking Llama2 Inference Serving on AMD EPYC Systems

The latency, throughput, and scaling testing described in this user guide uses EchoSwift\*.

## 5.1 Explaining Benchmarking

The objective of the LLM-Inference-Bench tools measures the latency of each request in milliseconds per token, Time Taken for the First Token (TTFT), and Throughput measured in the number of tokens per second of request sent to a served LLM. These metrics are captured using varying input tokens (query length), output tokens (response length), and simulated parallel users.



Figure 5-1: EchoSwift architecture

You can use this LLM inference benchmark to:

- **Measure performance:** Evaluate model latency, throughput, and resource (CPU-Core/Memory) utilization under varying workloads and configurations.
- **Identify bottlenecks:** Analyze performance metrics to pinpoint potential bottlenecks in different deployment scenarios impacting model efficiency and scalability.
- **Optimize deployment:** Utilize benchmarking data to fine-tune model deployment for improved resource utilization and cost-effectiveness.
- Validate scalability: Assess the model's ability to handle increasing workloads and maintain performance with a growing number of concurrent requests.

## AMD Nutanix<sup>™</sup> AOS/AHV Al-in-a-Box Solutions Build Guide

Metrics captured:

- Number of input tokens
- Number of output tokens
- Throughput (tokens/second)
- End-to-end latency (ms)
- Token latency (ms/tokens)
- TTFT for streaming
- CPU and memory utilization while running the load test.



Figure 5-2: Latency, throughput, and TTFT calculations

Performance metrics calculations:

- **Throughput:** Measures the average number of tokens that can be generated per second: Throughput (tokens/Second) = (Output Token Length) / Total Time (s)
- **Token latency:** Measures the average time it takes to generate one token: Latency (ms/token) = Total Time (s) \* 1000 / Output Token Length
- **TTFT:** Indicates the responsiveness of the model. Time from start of request to generation of first token.



Locust is an open-source tool that simulates multiuser requests on systems under test. It is built in python and excels at user behavior simulation. It is both flexible and easy to use, and its distributed testing capabilities allow simulating thousands of concurrent users. These features enable comprehensive system performance evaluation under realistic scenarios. The following prerequisites must be met in order to run the benchmark:

Item	<b>Configuration/Version</b>
# of cores	16
RAM	20 GB
datasets	2.13.1
locust	2.18.4
docker	5.0.3
kserve	0.10.0
torch	1.13.0
transformers	4.31.0
huggingface_hub	0.20.0

Table 5-1: Locust hardware and software requirements

#### 5.2.1 Benchmark Client

The benchmark script takes the required variables and starts sending requests to the deployed model endpoint. Locust creates virtual users and simulates production-level load testing. The benchmark will run against the above generation endpoint, which can be any inference server endpoint, such as a TGI endpoint hosting a model (http://localhost:8080/generate\_stream).

Define the configurations inside the locust.sh shell script before starting the load test. Here are some sample configurations:

- Number of parallel users [1, 3, 10].
- Varying input tokens [32, 64, 128]
- Varying output tokens [32, 64, 128]
- \$ ./llm\_inference\_benchmark.sh<output\_dir <generation\_endpoint>

## 



# Linear Llama2 7B "CPU-Only" Scaling

Deploying LLMs for real-time inference tasks often requires efficient resource management to handle varying workloads. This chapter discusses linear scaling of CPU-only LLM serving on Llama2 7B utilizing an infrastructure built on NKE. The primary objective is to showcase how the CPU resources allocated to Llama2 7B serving pods scale linearly with increasing inference requests.



Figure 6-1: Deployment options

Figure 6-1 shows a Nutanix LLM deployment using Kind deployment. The Deployment path offers a Horizontal Pod Autoscaler for scaling container instances and a Service Mesh for managing microservices. Autoscaling allows the number of serving replicas to adjust dynamically based on CPU and memory utilization, thereby ensuring optimal performance and resource utilization. However, the Nutanix cluster only supports Persistent volume with the access mode Read-Write-Once; scaling multiple replicas is thus only possible if the node has enough resources for all of the replicas, to avoid pod mount/volume issues. Please see <u>Persistent Volume Claims</u>\* for more information.