

# 计算加速套件 ( TACO )

## 产品文档



腾讯云TCE

## 目录

- 计算加速套件 (TACO) ..... 3
  - 产品简介 ..... 3
    - 产品概述 ..... 3
    - 产品优势 ..... 4
    - 应用场景 ..... 7
    - 产品架构 ..... 9
  - 功能特性 ..... 10
    - TACO Train AI训练加速 ..... 10
    - TACO LLM 推理加速引擎 ..... 14
      - 主要特性 ..... 14
      - 支持模型 ..... 20
  - 使用建议 ..... 24
  - 用户手册 ..... 25
    - TACO Train AI 训练加速引擎 ..... 25
      - 使用 TACO Train 部署分布式集群 ..... 25
        - 在 CVM 上部署 TensorFlow 分布式训练集群 ..... 25
        - 在裸金属服务器上部署 TensorFlow 分布式训练集群 ..... 33
        - 在 CVM 上部署 PyTorch 分布式训练集群 ..... 42
        - 在裸金属服务器上部署 PyTorch 分布式训练集群 ..... 51
    - 组件配置和使用 ..... 59
      - TCCL 使用说明 ..... 59
      - 配置 HARP 分布式训练环境 ..... 67
      - 配置容器 SSH 免密访问 ..... 69
      - 容器安装用户态 RDMA 驱动 ..... 71
  - TACO LLM 推理加速引擎 ..... 74
    - TACO-LLM 部署 ..... 74
    - TACO LLM 安装 ..... 79
    - TACO LLM 使用 ..... 80
      - 离线模式 ..... 80
      - 在线模式 ..... 82
      - 基础配置 ..... 84
        - Lookahead Cache ..... 86
        - Auto Prefix Caching ..... 93
        - 量化 ..... 96
      - CPU 辅助加速 ..... 101
      - 长序列优化 ..... 107
  - TACO LLM API ..... 111
    - Offline API ..... 111
    - Online API ..... 115
    - Sampling API ..... 123
- TACO LLM 性能 ..... 126

# 产品简介

## 产品概述

计算加速套件 TACO Kit ( Accelerated Computing Optimization Kit ) 是一种异构计算加速软件服务，搭配软硬件协同优化组件和硬件厂商特有优化方案，支持物理机、云服务器、容器等产品的计算加速、图形渲染、视频转码各个应用场景，帮助用户实现全方位全场景的降本增效。

# 产品优势

TACO Train AI 训练加速引擎和 TACO Infer AI 推理加速引擎通过软硬件协同优化，屏蔽底层硬件差异，适配 CPU、GPU、NPU 等不同加速硬件，降低用户使用计算资源的学习成本的同时提高计算效能。

## TACO Train AI 训练加速引擎优势

### 多层次深度优化加速

提供从自底向上的网络通信、分布式策略及训练框架等多层级的优化加速组件，用户可以根据需要选择适配。

### 支持无侵入式业务迁移

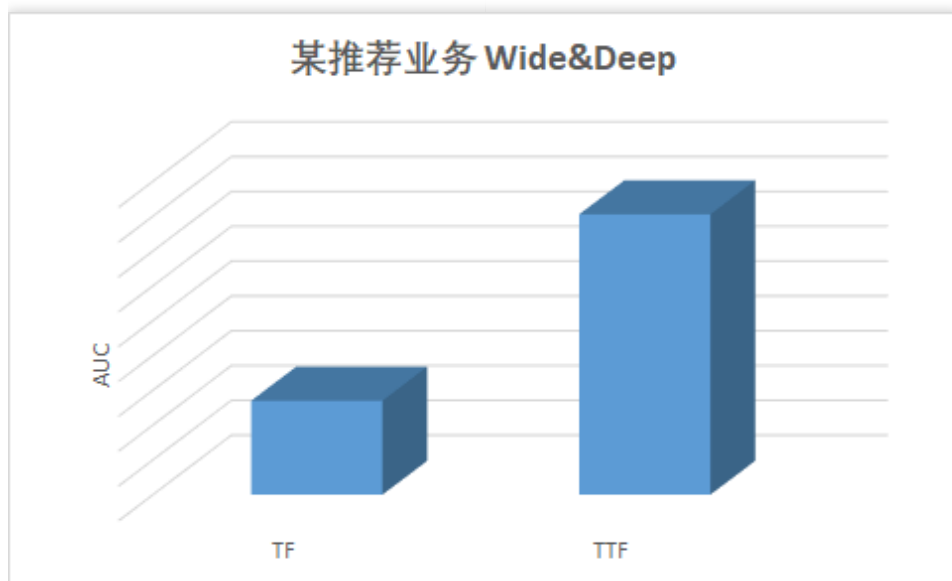
HARP、LightCC 等优化技术支持插件式集成，无需业务代码改动，即可加速分布式训练业务。

### 灵活扩展分布式训练场景

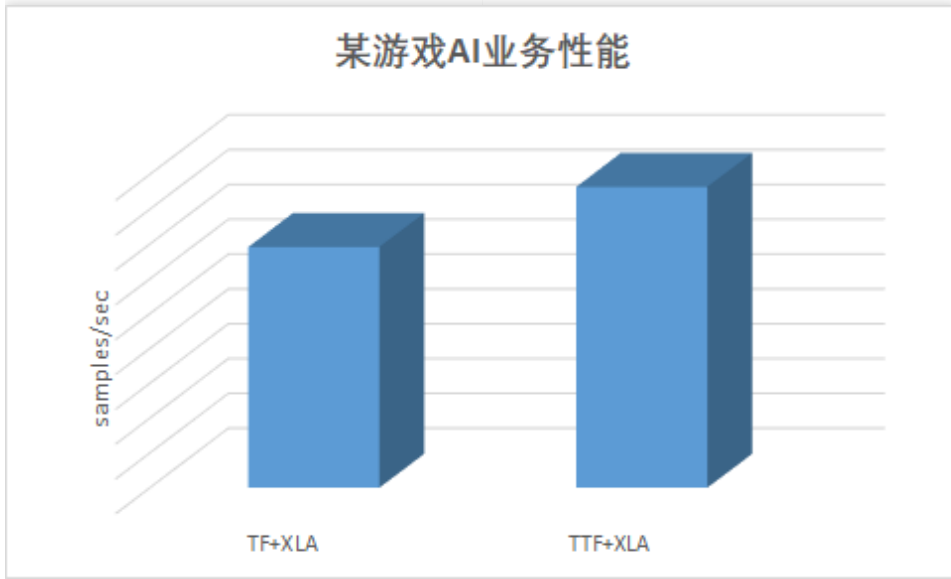
支持大规模多机多卡分布式训练场景，提高加速比和模型迭代效率。

### 训练性能提升数据

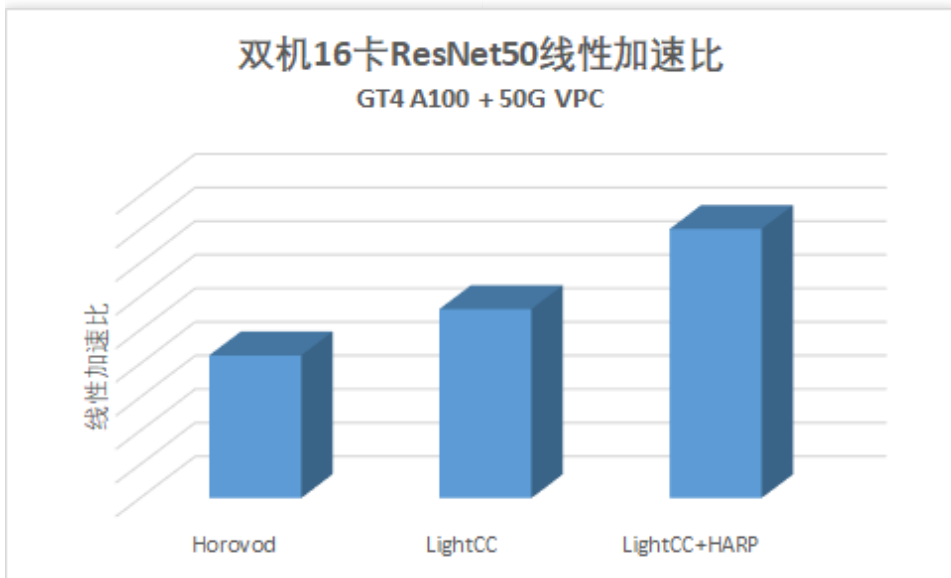
Tensorflow (以下简称 TTF) 动态 Embedding 在某推荐业务上对 AUC 的提升效果：



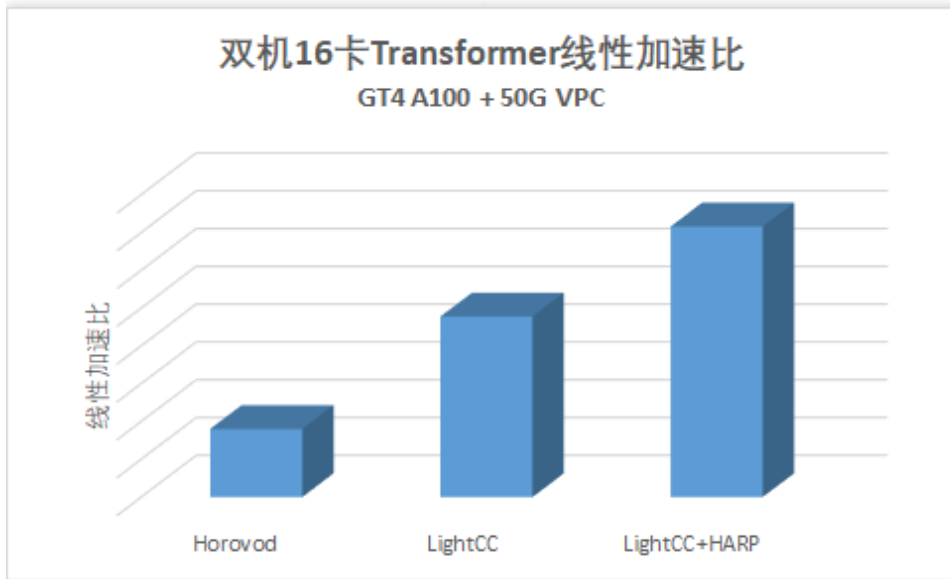
TTF XLA 在某游戏业务上的性能加速效果：



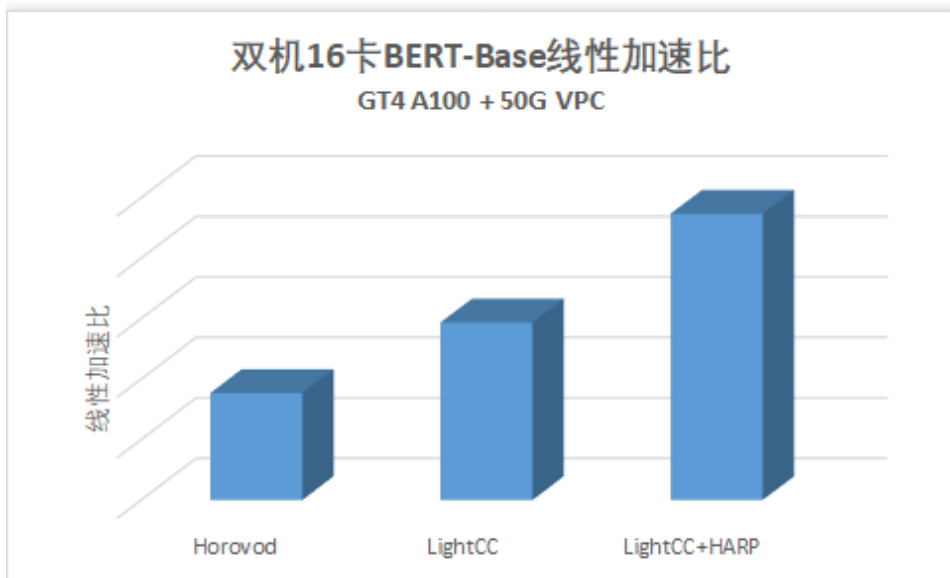
在50G VPC 环境下，ResNet50的多机训练加速效果：



在50G VPC 环境下，Transformer 的多机训练加速效果：



在50G VPC 环境下，BERT-Base 的多机训练加速效果：



## TACO Infer AI 推理加速引擎优势

### 部署简洁

TACO Infer 仅有一行简洁的优化接口，不会改变用户的模型格式。用户可以保持其一贯的使用和部署习惯，并提供插件式的第三方开发接口，支持适配不同业务场景。

### 软硬件兼容

支持多种框架模型和多种加速硬件，可运行在虚拟机、物理机、容器等各种环境。

### 一站式解决推理部署相关优化依赖

集成硬件厂商的定向开源的加速方案，整合先进的编译优化、图优化和算子优化技术。

# 应用场景

TACO Train AI 训练加速引擎和 TACO Infer AI 推理加速引擎目前支持但不限于以下场景：

计算机视觉：例如 ResNet、MobileNet、Inception、ViT 等。

推荐系统：例如 Wide&Deep、DeepFM 等。

NLP 服务：例如 BERT、Transformer、ASR 等。

TACO-LLM 适用于大语言模型的推理加速业务，可满足多种业务场景下推理提效的需求。以下是一些典型业务场景：

## 客户服务

业务场景	场景解释
智能客服	用于回答客户咨询，提供24小时服务。
问答系统	自动识别客户意图并作出响应。
情感分析	分析客户反馈和评论，识别客户情绪，帮助企业改进服务。

## 内容创作与编辑

业务场景	场景解释
文本润色	对用户提供的文本进行润色加工，生成质量更高的文本内容。
文本摘要	提取长篇文章的要点，生成摘要。
文本写作	生成新闻报道、文章、博客等内容。

## 翻译与本地化

业务场景	场景解释
机器翻译	将一种语言的文本翻译成另一种语言。
内容本地化	根据不同地区的文化和习惯调整内容。

## 编程开发

业务场景	场景解释
代码助手	根据自然语言描述或者已有代码片段自动生成代码，辅助开发者编程。
代码审查	自动检查代码中的错误和潜在问题。

## 教育培训

业务场景	场景解释
学习助手	根据用户的需求提供定制化的学习材料和解答。
RAG 知识引擎	结合外部知识库和大语言模型的能力，构建出强大的知识引擎。

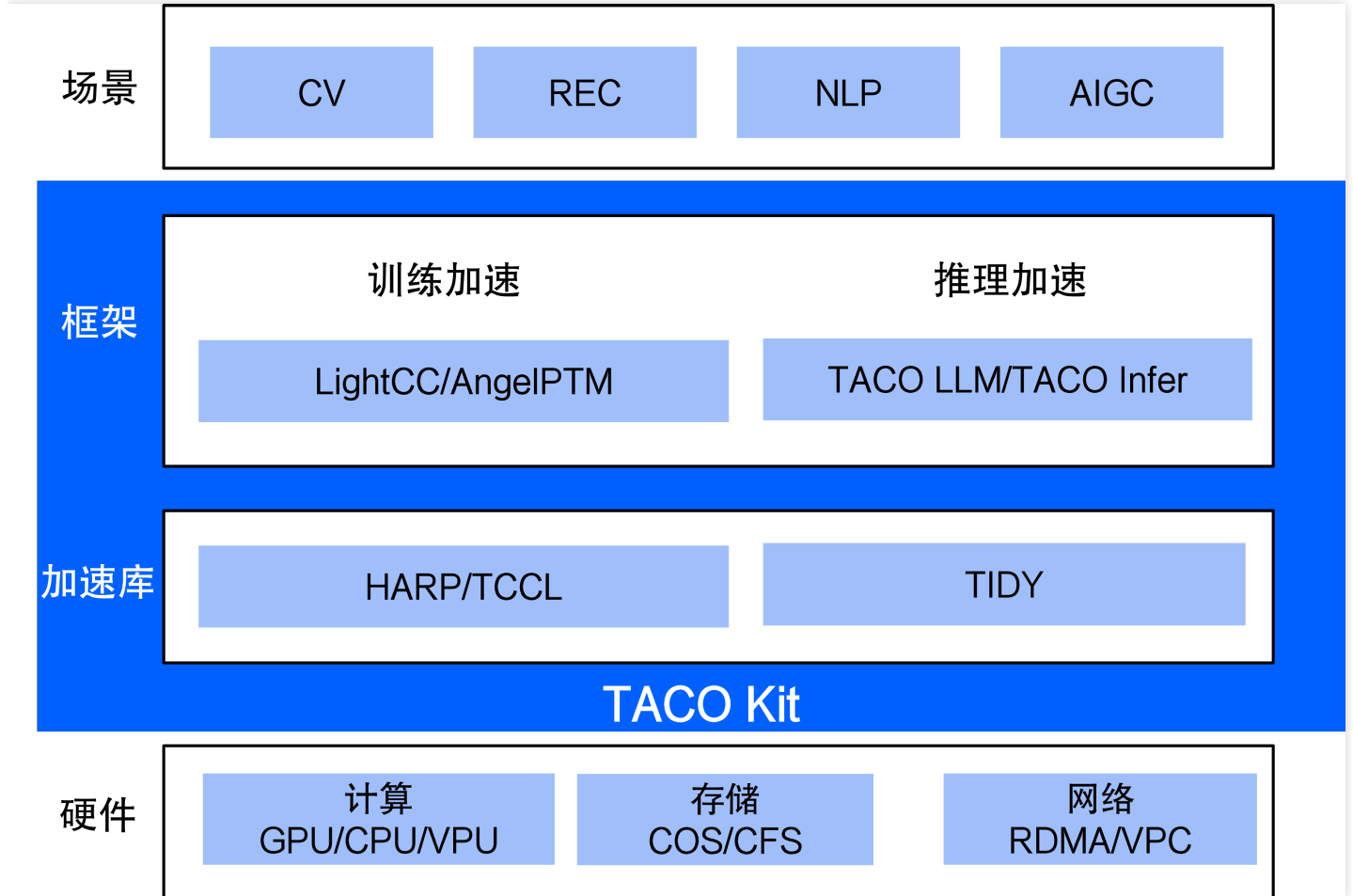
## LLM 训练辅助

业务场景	场景解释
预训练数据生成	利用大语言模型辅助生成大量预训练数据，供后续训练大语言模型使用。

TACO-LLM 在对时延敏感的在线服务场景(如智能客服和问答系统)和要求高吞吐的离线服务场景(如预训练数据生成和文本摘要)都提供了极具竞争力的性能加速方案，可帮助您最大效率地利用算力资源实现高吞吐，低延时，大幅降低单任务平均算力成本，获得最佳成本优化方案。

# 产品架构

TACO Kit 支持丰富的 AI 业务场景，深耕训练框架优化、分布式框架优化、网络通信优化、推理性能优化等关键技术，致力于打造整套的 AI 加速方案。其架构示意图如下所示：



# 功能特性

## TACO Train AI训练加速

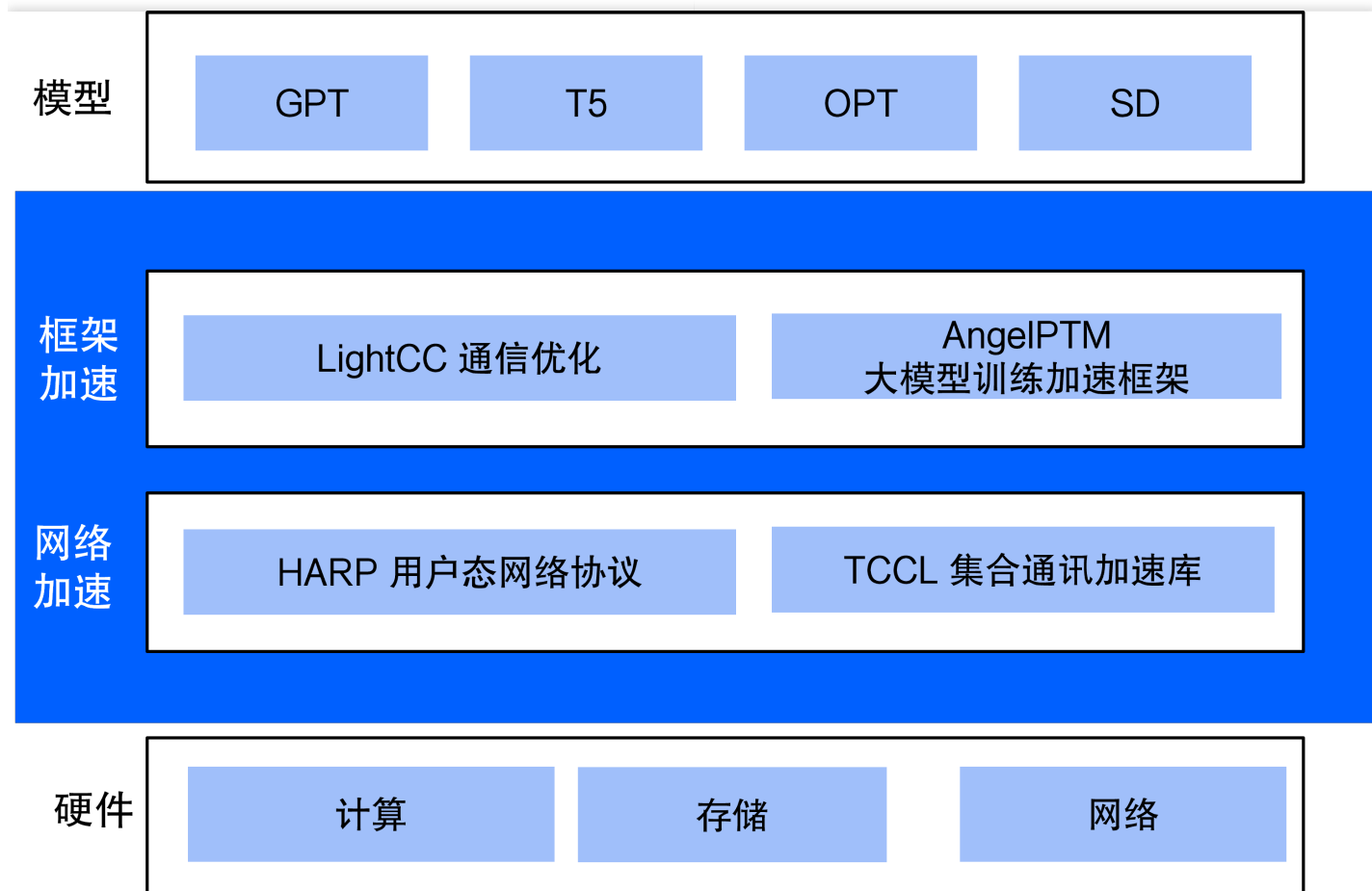
### 背景信息

近几年随着 AI 模型参数的倍增及训练数据的日益增长，用户对模型迭代效率的需求也随之增长，单个 GPU 的算力和显存资源已无法满足大部分业务场景，使用单机多卡或多机多卡训练已成为趋势。单机多卡训练场景的参数同步借助 NVIDIA NVLINK 技术，基本可以获得较高的线性扩展比，但多机多卡训练场景严重依赖多机之间的网络互联技术。网卡厂商提供了高速互联技术 Infiniband 或 RoCE，虽使多机通信效率大幅提升，但成本也大幅增加。如何在普通甚至低速网络环境下提升分布式系统的训练效率，已成为用户关注的焦点。

### TACO Train 简介

目前业内已有一些成熟的分布式训练加速技术，例如多级通信、多流通信、梯度融合、压缩通信等，TACO Train 也引入了类似的加速技术。同时，TACO Train 推出了自定义的用户态协议栈 HARP，有效解决普通网络环境下的多机网络通信问题。

TACO Train 基于 IaaS 资源推出的 AI 训练加速引擎，为用户提供开箱即用的 AI 训练套件。TACO Train 基于丰富的 AI 业务场景，提供自底向上的网络通信、分布式策略及训练框架等多层级的优化，是一套全生态的训练加速方案。



目前 TACO Train 提供了4个训练加速组件：

- AngelPTM：大模型预训练框架，支持 GPT/T5/BERT 等大模型。
- TCCL：针对星脉网络架构的高性能定制加速通信库，为 AI 大模型训练提供更高效率的网络通信性能。
- HARP：自研用户态网络协议栈，提高普通网络环境的多机通信效率。
- LightCC：分布式训练框架，支持 Horovod，Ray 或者 PyTorch DDP 的分布式训练加速。

## AngelPTM

AngelPTM 是基于 DeepSpeed 和 Megatron 深度定制开发的大模型训练框架，支持 NLP/ 多模态 /AIGC 等多类预训练任务。由于大模型的参数规模巨大，对硬件存储资源提出了挑战，AngelPTM 支持以更少的资源和更快的速度训练大模型，兼容社区方案 API，支持业务快速接入，

- ZeRO Cache 策略：基于 ZeRO 策略把内存作为二级存储 offload 参数、梯度、优化器状态到 CPU 内存（也支持 SSD 作为第三级存储）。ZeRO-Cache 为了最大化利用内存和显存进行模型状态的缓存，引入了显存内存统一存储视角，将存储容量的上界由内存扩容到内存+显存总和。同时将多流异步化做到了极致，在 GPU 计算的同时进行数据 IO 和 NCCL 通信，使用异构流水线均衡设备间的负载，最大化提升整个系统的吞吐。ZeRO-Cache 将 GPU 显存、CPU 内存统一视角管理，击破了异构存储的壁垒，减少了冗余存储和内存碎片，增加了内存的利用率，极大扩充了模型存储可用空间。

- 自动流水并行：3D 并行是 Megatron 的原始能力，但是 Megatron 的流水并行只支持 block 级别（Transformer Layer），且需要依靠专家经验人工指定 stage 切分以确保 stage 间负载均衡；自动流水并行可以做到 op 级别，且自动 profile op 性能，自动搜索负载均衡的切分方案后执行流水并行训练。
- 高性能 MoE 组件：基于拓扑感知的 AlltoAll 通信加速专家并行，提供高性能定制算子，支持计算通信流水提高迭代效率。

## TCCL

TCCL 是一款针对星脉网络架构的高性能定制加速通信库。主要功能是依托星脉网络硬件架构，为 AI 大模型训练提供更高的网络通信性能，同时具备网络故障快速感知与自愈的智能运维能力。TCCL 基于开源的 NCCL 代码做了扩展优化，完全兼容 NCCL 的功能与使用方法。TCCL 目前支持主要特性包括：

- 双网口动态聚合优化，发挥 bonding 设备的性能极限。
- 全局 Hash 路由（Global Hash Routing），负载均衡，避免拥塞。
- 拓扑亲和性流量调度，最小化流量绕行。

## HARP

随着网络硬件技术的发展，网络带宽从 10Gbps 增长到 100Gbps 甚至更高，在数据中心大量部署使用。但目前普遍使用的内核网络协议栈存在着一些必要的开销，使其不能很好地利用高速网络设备。为解决该问题，自研了用户态网络协议栈 HARP，可以以 Plug-in 的方式集成到 NCCL 中，无需任何业务改动，加速云上分布式训练性能。在 VPC 的环境下，相比传统的内核协议栈，HARP 提供了以下的能力：

- 支持全链路内存零拷贝，HARP 协议栈提供特定的 buffer 给应用，使应用的数据经过 HARP 协议栈处理后由网卡直接进行收发，消除内核协议栈中耗时及占用 CPU 较高的多次内存拷贝操作。
- 支持协议栈多实例隔离，即应用可以在多个 CPU core 上创建特定协议栈实例处理网络报文，每个实例间相互隔离，保证性能线性增长。
- 数据平面无锁设计，HARP 协议栈内部保证网络 session 的数据仅在创建该 session 的 CPU core 上，使用特定的协议栈实例处理。减少了内核中同步锁的开销，也降低了 CPU 的 Cache Miss 率，大幅提升网络数据的处理性能。

## LightCC

LightCC 对社区分布式方案的通信策略进行了深度定制优化，完全兼容 Horovod，Ray 或者 PyTorch DDP API，支持业务快速接入。主要包括的优化能力如下：

- 2D AllReduce 充分利用通信带宽。
- 高效的梯度融合方式。

- TOPK/FP16 压缩通信，降低通信量，提高传输效率。
- 简单高效地管理和分配分布式训练任务，具有较强的扩展性和可靠性。

## TTF

TensorFlow 是深度学习领域中应用最广泛的开源框架之一，但在很多业务场景下，开源 TensorFlow 有其特定的限制。为了解决实际业务中遇到的问题，Tensorflow (以下简称 TTF) 提供了以下能力：

- 相比原始的静态 Embedding，高维稀疏动态 Embedding 帮助用户在不需要重新训练的条件下，动态添加和删除特征，按需使用内存，避免 Hash 冲突，同时保留原始 TF 的 API 设计风格。
- 混合精度在原有实现的基础上增加了调整精度的策略，根据 loss 的状态自动在全精度和半精度之间切换，避免精度损失。
- 针对特定业务场景的 XLA，Grapppler 图优化，以及自适应编译框架解决冗余编译的问题。
- 开源 TF 1版本不再提供对 Ampere GPU 的支持，但考虑到较多用户仍在使用 TF 1.15版本的问题，TTF 添加了对 CUDA 11的支持，让用户可以使用 A100来进行模型训练。

# TACO LLM 推理加速引擎

## 主要特性

### LLM 服务部署的挑战

与传统 AI 场景不同，大语言模型服务关注多个维度的性能指标，它们最终对应不同层次的用户体验和服务质量。这些指标大体概括起来分为4个：

- 首字延迟 (L1)

定义为 LLM 服务处理完 prompt 输出第一个 token 的延时，决定了用户从输入请求到获得响应的时间。对于实时的在线应用，低延迟很重要。但对偏离线的应用则该指标没有那么重要。该指标的好坏通常取决于推理引擎处理 prompt 并生成首字的时间。

- 解码延迟 (L2)

定义为每个用户请求生成后续输出的平均响应时间，也是用户直观体验模型执行快慢的时间。假设执行速度是平均每 Token 100毫秒，则用户平均每秒可以得到10个 Tokens 的输出，每分钟 450 左右英文词。

- 请求延迟 (L3)

定义为对给定用户产生完整响应的延迟。计算规则： $L3 = L1 + L2 \times \text{生成的 Token 数}$ 。

- 吞吐

定义为推理服务器面对全部用户和他们请求的流量时每秒可以生成的 Token 数量。

部分推理引擎只关注或对上述某个指标有较好效果。而 TACO-LLM 均衡关注上述全部指标，并对各指标的实际部署效果均实现了全流程的优化。

LLM 部署的挑战来源于几个方面：

- 当前主流的 decoder-only 模型都具备自回归解码属性。模型生成输出是一个串行的计算过程。下一个输出依赖上一个输出。因此很难发挥出 GPU 或其他加速硬件的并行加速能力。同时，较低的 Arithmetic Intensity 对显存带宽的利用也提出了挑战。
- 大模型的大对显存容量提出了最直接的挑战。
- 传统的 Transformer 推理框架将 KV-Cache 按 batchsize 和 sequence length 维度组织数据，这会导致两个潜在的性能陷阱：

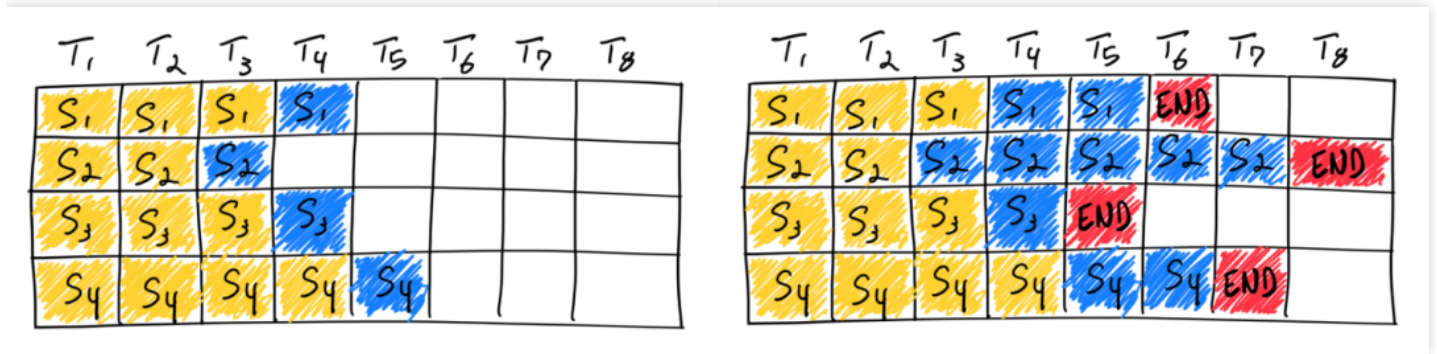
- 请求的输出有长有短。这种数据组织方式需要等一个 Batch 中最长输出长度的请求计算完才能完成整个 Batch 的计算。在此之前，新的请求无法开始计算。而已计算完的请求，只能进行无效计算，消耗有效算力。同时，这种方式管理显存效率较低，无法做到“实用实销”，且随着不同请求计算过程中的显存使用，会造成显存碎片，进一步加剧资源瓶颈。
- 目前很多优秀的 attention 加速技术实际上是按上述数据 layout 来实现的。例如 flash-attention、flash-decoding 等。这意味着更高效的重新设计，例如下文中提到的 Paged Attention 技术将无法直接享受到社区优化红利，仍给我们留出了进一步的优化空间。

面向生产的 LLM 推理引擎有效的应对了上述挑战，全方位解决问题。

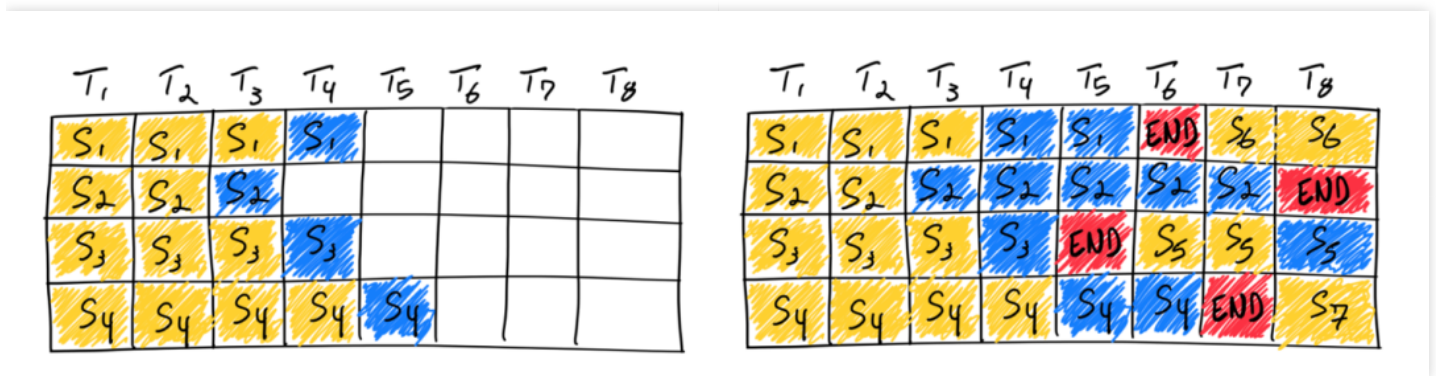
# Continuous Batching

传统的 Batching 方式被称为 Static Batching。如上文所述，Static Batching 方式需要等一个 Batch 中最长输出长度的请求完成计算，整个 Batch 才完成返回，新的请求才能重新 Batch 并开始计算。因此，Static Batching 方式在其他请求计算完成，等待最长输出请求计算的过程中，严重浪费了硬件算力。TACO-LLM 通过 Continuous Batching 的方式来解决这个问题。Continuous Batching 无需等待 Batch 中所有请求都完成计算，而是一旦有请求完成计算，即可以加入新的请求，实现迭代级别的调度，提高计算效率。从而实现较高的 GPU 计算利用率。

## Static Batching



## Continuous Batching



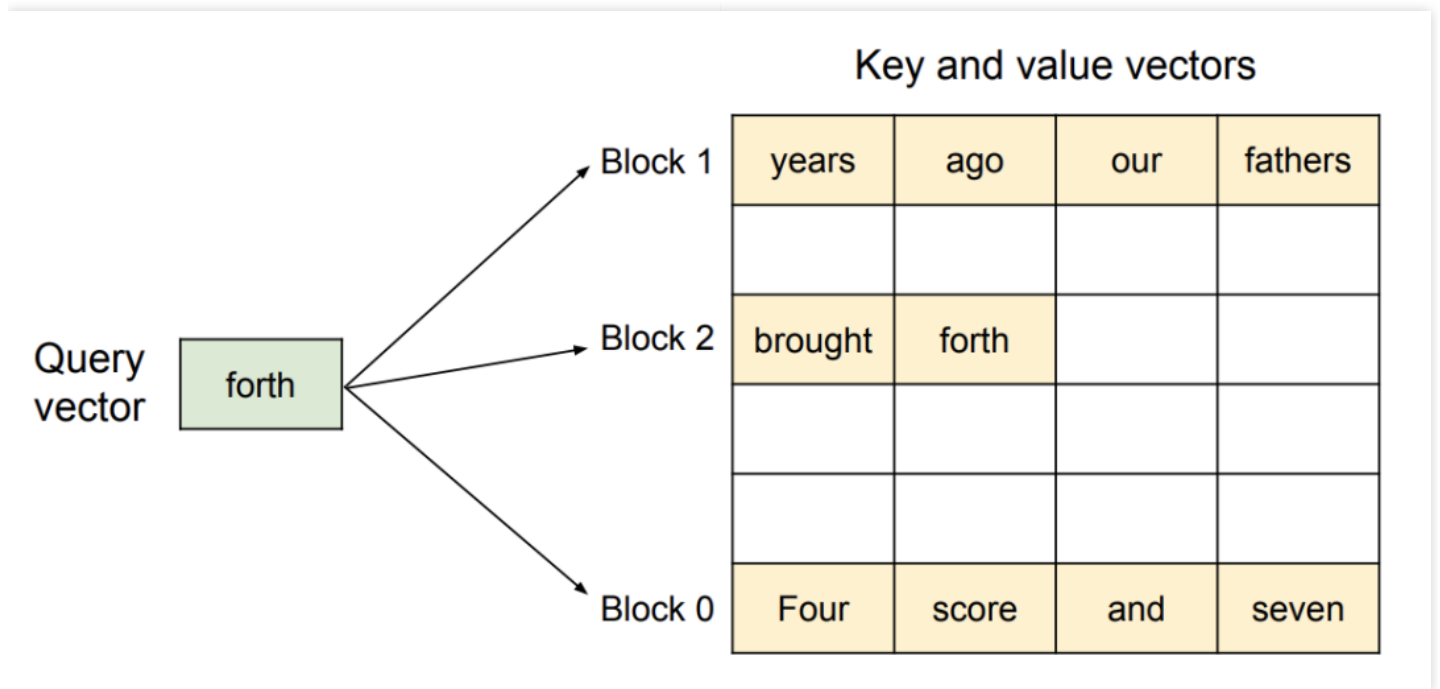
# Paged Attention

大模型推理计算性能优化一个常用的方式是 KV-Cache 技术。Transformer 层的 attention 组件计算当前 Token value 值时，需要依赖之前 Token 序列的 Key 和 Value 值。KV-Cache 通过存储之前 Token 序列的 Key 和 Value 的值，避免后续计算中，重复计算 Key 和 Value 值，提高整体计算性能，是一种以空间换时间的优化策略。

传统的 KV-Cache 实现机制是在显存中提前预留一块连续的存储空间来存储 Key 和 Value 值。但是，随着存储资源的分配和释放，显存中会存在很多“碎片”。某些情况下，虽然剩余的显存总量大于 KV-Cache 所需，但是由于不存在一块连续的存储空间可以满足 KV-Cache，计算也无法进行。

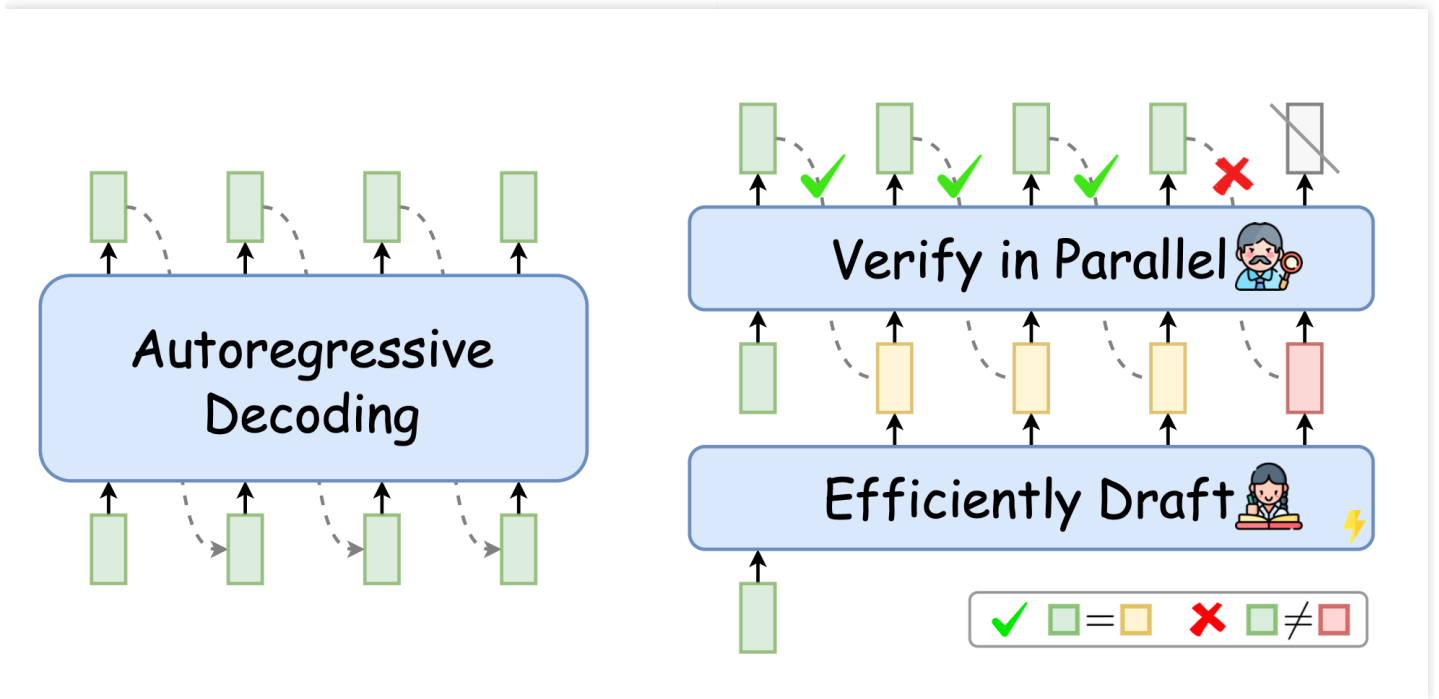
Paged Attention 是一种新的 KV-Cache 实现方式，它从传统操作系统的概念中获得灵感，例如分页和虚拟内存，允许 KV-Cache 通过分配固定大小的“页”或“块”在物理非连续内存上实现逻辑连续。然后可以将注意力机制重写为在块

对齐的输入上运行，从而允许在非连续的内存范围内执行注意力计算。TACO-LLM 通过 Paged Attention 技术，实现了较高的显存利用效率。



## 投机采样

大语言模型的自回归解码属性要求每次生成新的 Token，都需要依赖所有已解码的 Token，且需要重新加载模型全部权重进行串行解码。这种计算方式无法充分利用 GPU 的算力，计算效率不高，解码成本高昂。而 TACO-LLM 通过投机采样的方式，从根本上解决了计算访存比的问题。通过引入一个高效的 Draft 辅助解码，可以让真正部署的大模型实现“并行”解码，从而大幅提高解码效率。我们称之为 **Speculative Sampling (SpS)** 技术。

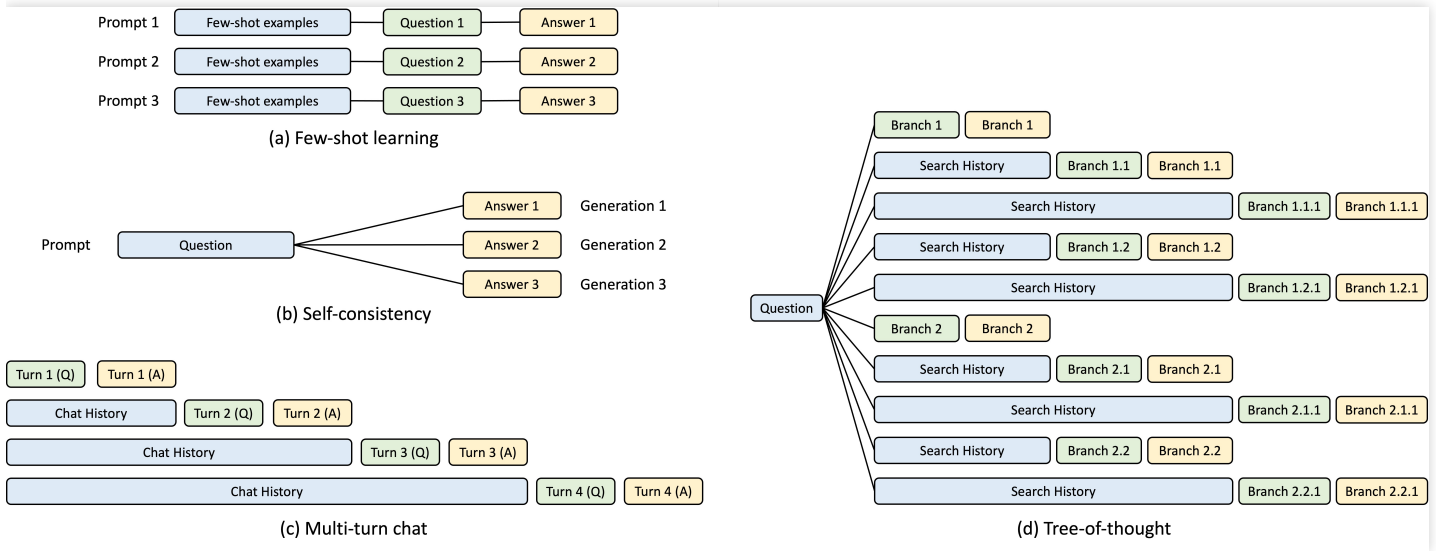


TACO-LLM 支持多种投机采样技术，如基于上下文的 Lookahead-Cache 和基于模型的 Eagle/Medusa 等，在多种业务场景中，均能实现高效的 LLM 推理计算。关于 Lookahead-Cache 的详情和使用方式请参见 Lookahead Cache。

## Auto Prefix Caching

LLM 推理计算主要分为两个过程：**Prefill 阶段 (Prompt 计算)** 和 **Decode 阶段**。这两个阶段的计算特性存在不同，Prefill 阶段是计算受限的，而 Decode 阶段是访存受限的。为了避免重复计算，Prefill 阶段主要作用就是给 Decode 阶段准备 KV Cache。但这些 KV Cache 通常只是为单条推理请求服务的，当请求结束，对应的 KV-Cache 就会清除。

KV Cache 能不能跨请求复用？在某些 LLM 业务场景下，多次请求的 Prompt 可能会共享同一个前缀 (Prefix)，比如少量样本学习，多轮对话等。在这些情况下，很多请求 Prompt 的前缀的 KV Cache 计算的结果是相同的，可以被缓存起来，给之后的请求复用。TACO-LLM 的 Auto Prefix Cache 技术可以针对这种场景进行优化，使得具有相同 Prompt 前缀的 KV-Cache 可以跨请求复用，降低计算开销，提升推理计算性能。详情和使用方式请参见 Auto Prefix Cache。



## 量化

随着以 Transformer 为基石的大语言模型规模的快速增大，LLM 的推理部署对 GPU 显存和算力的需求激增。如何减少大语言模型部署的 GPU 显存需求，提高计算效率，降低推理部署成本就显得愈发重要。模型量化是解决这个问题非常重要的一种手段。业界提出了多种适用于 LLM 的量化算法，如 GPTQ、AWQ、FP8 等，在保证模型精度损失满足业务需求的情况下，显著提升 LLM 推理计算的效能，降低部署成本。TACO-LLM 对业界主要的 LLM 量化算法均进行了适配支持，详情和使用方式请参见 [量化](#)。

## CPU 辅助加速

传统的投机采样使用 GPU 作为 Draft Model 的计算资源，而 GPU 的成本高昂。为了进一步降低计算成本，TACO 团队与 Intel 团队合作，基于 AMX 指令集对 CPU 上的矩阵乘法做了优化，使得使用 CPU 作为 Draft Model 成为可能，从而在进行推理加速的同时，显著降低了推理成本。详情和使用方式请参见 [CPU 辅助加速](#)。

## 长序列并行

在 LLM 大模型推理中，长序列场景应用越来越广泛。目前业界对长序列的优化主要有 KV-Cache 量化、稀疏化等，这些都对模型精度有一定的影响。TACO LLM 的长序列并行方案，可以在长序列场景进行精度无损的加速。在长序列推理场景，LLM 推理的首字延迟较高。针对该问题，序列并行在 Prefill 阶段采用了 Ring Attention 类的方案，可以通过扩展机器，降低首字延迟。而对于推理阶段，我们可以使用 Lookahead 等投机采样技术加速。详情和使用方式请参见 [长序列优化](#)。

## 参考文献

- [1] [How continuous batching enables 23x throughput in LLM inference while reducing p50 latency](#)
- [2] [Efficient Memory Management for Large Language Model Serving with PagedAttention](#)
- [3] [Unlocking Efficiency in Large Language Model Inference: A Comprehensive Survey of Speculative Decoding](#)
- [4] [Fast and Expressive LLM Inference with RadixAttention and SGLang](#)

# 支持模型

TACO-LLM 支持 Huggingface 模型格式的多种生成式 Transformer 模型。下面列出了 TACO-LLM 目前支持的模型架构和对应的常用模型。

## Decoder-only 语言模型

Architecture	Models	Example HuggingFace Models	LoRA
BaiChuanForCausalLM	Baichuan & Baichuan2	baichuan-inc/Baichuan2-13B-Chat, baichuan-inc/Baichuan-7B, etc.	Supported
BloomForCausalLM	BLOOM, BLOOMZ, BLOOMChat	bigscience/bloom, bigscience/bloomz, etc.	-
ChatGLMModel	ChatGLM	THUDM/chatglm2-6b, THUDM/chatglm3-6b, etc.	Supported
FalconForCausalLM	Falcon	tiiuae/falcon-7b, tiiuae/falcon-40b, tiiuae/falcon-rw-7b, etc.	-
GemmaForCausalLM	Gemma	google/gemma-2b, google/gemma-7b, etc.	Supported
Gemma2ForCausalLM	Gemma2	google/gemma-2-9b, google/gemma-2-27b, etc.	Supported
GPT2LMHeadModel	GPT-2	gpt2, gpt2-xl, etc.	-
GPTBigCodeForCausalLM	StarCoder, SantaCoder, WizardCoder	bigcode/starcoder, bigcode/gpt_bigcode-santacoder, WizardLM/WizardCoder-15B-V1.0, etc.	Supported
GPTJForCausalLM	GPT-J	EleutherAI/gpt-j-6b, nomic-ai/gpt4all-j, etc.	-
GPTNeoXForCausalLM	GPT-NeoX, Pythia, OpenAssistant, Dolly V2, StableLM	EleutherAI/gpt-neox-20b, EleutherAI/pythia-12b, OpenAssistant/oasst-sft-4-pythia-12b-epoch-3.5, databricks/dolly-v2-12b, stabilityai/stablelm-tuned-alpha-7b, etc.	-

Architecture	Models	Example HuggingFace Models	LoRA
InternLMForCausalLM	InternLM	internlm/internlm-7b, internlm/internlm-chat-7b, etc.	Supported
InternLM2ForCausalLM	InternLM2	internlm/internlm2-7b, internlm/internlm2-chat-7b, etc.	-
LlamaForCausalLM	Llama 3.1, Llama 3, Llama 2, LLaMA, Yi	meta-llama/Meta-Llama-3.1-405B-Instruct, meta-llama/Meta-Llama-3.1-70B, meta-llama/Meta-Llama-3-70B-Instruct, meta-llama/Llama-2-70b-hf, 01-ai/Yi-34B, etc.	Supported
MistralForCausalLM	Mistral, Mistral-Instruct	mistralai/Mistral-7B-v0.1, mistralai/Mistral-7B-Instruct-v0.1, etc.	Supported
MixtralForCausalLM	Mixtral-8x7B, Mixtral-8x7B-Instruct	mistralai/Mixtral-8x7B-v0.1, mistralai/Mixtral-8x7B-Instruct-v0.1, mistral-community/Mixtral-8x22B-v0.1, etc.	Supported
NemotronForCausalLM	Nemotron-3, Nemotron-4, Minitron	nvidia/Minitron-8B-Base, mgoin/Nemotron-4-340B-Base-hf-FP8, etc.	Supported
OPTForCausalLM	OPT, OPT-IML	facebook/opt-66b, facebook/opt-impl-max-30b, etc.	
PhiForCausalLM	Phi	microsoft/phi-1_5, microsoft/phi-2, etc.	Supported
Phi3ForCausalLM	Phi-3	microsoft/Phi-3-mini-4k-instruct, microsoft/Phi-3-mini-128k-instruct, microsoft/Phi-3-medium-128k-instruct, etc.	-
Phi3SmallForCausalLM	Phi-3-Small	microsoft/Phi-3-small-8k-instruct, microsoft/Phi-3-small-128k-instruct, etc.	-
PhiMoEForCausalLM	Phi-3.5-MoE	microsoft/Phi-3.5-MoE-instruct, etc.	-
QWenLMHeadModel	Qwen	Qwen/Qwen-7B, Qwen/Qwen-7B-Chat, etc.	-
Qwen2ForCausalLM	Qwen2	Qwen/Qwen2-beta-7B, Qwen/Qwen2-beta-7B-Chat, etc.	Supported

Architecture	Models	Example HuggingFace Models	LoRA
Qwen2MoeForCausalLM	Qwen2MoE	Qwen/Qwen1.5-MoE-A2.7B, Qwen/Qwen1.5-MoE-A2.7B-Chat, etc.	-
StableLmForCausalLM	StableLM	stabilityai/stablelm-3b-4e1t/ , stabilityai/stablelm-base-alpha-7b-v2, etc.	-
Starcoder2ForCausalLM	Starcoder2	bigcode/starcoder2-3b, bigcode/starcoder2-7b, bigcode/starcoder2-15b, etc.	-
XverseForCausalLM	Xverse	xverse/XVERSE-7B-Chat, xverse/XVERSE-13B-Chat, xverse/XVERSE-65B-Chat, etc.	-

## 多模态语言模型

Architecture	Models	Modalities	Example HuggingFace Models	LoRA
InternVLChatModel	InternVL2	Image(E+)	OpenGVLab/InternVL2-4B, OpenGVLab/InternVL2-8B, etc.	-
LlavaForConditionalGeneration	LLaVA-1.5	Image(E+)	llava-hf/llava-1.5-7b-hf, llava-hf/llava-1.5-13b-hf, etc.	-
LlavaNextForConditionalGeneration	LLaVA-NeXT	Image(E+)	llava-hf/llava-v1.6-mistral-7b-hf, llava-hf/llava-v1.6-vicuna-7b-hf, etc.	-
LlavaNextVideoForConditionalGeneration	LLaVA-NeXT-Video	Video	llava-hf/LLaVA-NeXT-Video-7B-hf, etc. (see note)	-
PaliGemmaForConditionalGeneration	PaliGemma	Image(E)	google/paligemma-3b-	-

			pt-224, google/paligemma-3b-mix-224, etc.	
Phi3VForCausalLM	Phi-3-Vision, Phi-3.5-Vision	Image(E+)	microsoft/Phi-3-vision-128k-instruct, microsoft/Phi-3.5-vision-instruct etc.	-
PixtralForConditionalGeneration	Pixtral	Image(+)	mistralai/Pixtral-12B-2409	-
QWenLMHeadModel	Qwen-VL	Image(E+)	Qwen/Qwen-VL, Qwen/Qwen-VL-Chat, etc.	-
Qwen2VLForConditionalGeneration	Qwen2-VL (see note)	Image(+) / Video(+)	Qwen/Qwen2-VL-2B-Instruct, Qwen/Qwen2-VL-7B-Instruct, Qwen/Qwen2-VL-72B-Instruct, etc.	-

#### 说明：

- E: 表示 Pre-computed embeddings 可以作为多模态输入。
- +: 表示一个 prompt 可以插入多个多模态输入。

# 使用建议

## 软件使用建议

- TACO Infer AI 推理加速引擎硬件支持 Intel 和 AMD 系列的 CPU 优化，后续会支持更多优化目标硬件，请您关注产品动态。
- AI 框架已对 TensorFlow 1.14和1.15两个主要版本进行了全面的测试。如使用其他版本遇到任何问题，请联系工程师获取支持。
- 在使用 Keras+Horovod 训练过程中，建议使用 TCMalloc 进行内存优化，可以使内存在多个 epoch 之间基本保持不变。

# 用户手册

## TACO Train AI 训练加速引擎

### 使用 TACO Train 部署分布式集群

### 在 CVM 上部署 TensorFlow 分布式训练集群

## 操作场景

本文介绍如何基于云服务器 CVM 搭建 Tensorflow+Taco Train 分布式训练集群。

## 操作步骤

### 开通实例

开通实例，其中实例、存储及镜像请参考以下信息选择，其余配置请参考 [通过开通页创建实例](#) 按需选择。

- 实例：选择 对应的GPU实例。
- 系统盘：配置容量不小于50GB的云硬盘。您也可在创建实例后使用文件存储，详情参见 [在 Linux 客户端上使用 CFS 文件系统](#)。
- 镜像：建议选择公共镜像，您也可选择自定义镜像。
- 操作系统请使用 CentOS 8.0/CentOS 7.8/Ubuntu 20.04/Ubuntu 18.04。

### 配置实例环境

#### 验证 GPU 驱动

1. 登录实例。
2. 执行以下命令，验证 GPU 驱动是否安装成功。

```
nvidia-smi
```

查看输出结果是否为 GPU 状态：

- 是，代表 GPU 驱动安装成功。
- 否，请参考 [NVIDIA Driver Installation Quickstart Guide](#) 进行安装。

#### 配置 HARP 分布式训练环境

1. 参考 [配置 HARP 分布式训练环境](#)，配置所需环境。
2. 配置完成后，执行以下命令进行验证，若配置文件存在，则表示已配置成功。

```
ls /usr/local/tfabric/tools/config/ztcp*.conf
```

## 安装 docker 和 nvidia docker

1. 执行以下命令，安装 docker。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-docker.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 Docker 官方文档进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
+ sh -c 'yum install -y -q docker-ce-rootless-extras'
Package docker-ce-rootless-extras-20.10.16-3.el7.x86_64 already installed and latest version
=====

To run Docker as a non-privileged user, consider setting up the
Docker daemon in rootless mode for your user:

    dockerd-rootless-setuptool.sh install

Visit https://docs.docker.com/go/rootless/ to learn about rootless mode.

To run the Docker daemon as a fully privileged service, but granting non-root
users access, refer to https://docs.docker.com/go/daemon-access/

WARNING: Access to the remote API on a privileged Docker daemon is equivalent
to root access on the host. Refer to the 'Docker daemon attack surface'
documentation for details: https://docs.docker.com/go/attack-surface/
=====
```

2. 执行以下命令，安装 nvidia-docker2。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-nvidia-docker2.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 NVIDIA 官方文档 [Installation Guide & mdash](#) 进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
Running transaction
Installing : libnvidia-container1-1.9.0-1.x86_64
Installing : libnvidia-container-tools-1.9.0-1.x86_64
Installing : nvidia-container-toolkit-1.9.0-1.x86_64
Installing : nvidia-docker2-2.10.0-1.noarch
Verifying  : libnvidia-container-tools-1.9.0-1.x86_64
Verifying  : nvidia-container-toolkit-1.9.0-1.x86_64
Verifying  : nvidia-docker2-2.10.0-1.noarch
Verifying  : libnvidia-container1-1.9.0-1.x86_64

Installed:
nvidia-docker2.noarch 0:2.10.0-1

Dependency Installed:
libnvidia-container-tools.x86_64 0:1.9.0-1      libnvidia-container1.x86_64 0:1.9.0-1      nvidia-container-toolkit.x86_64 0:1.9.0-1

Complete!
```

## 下载 docker 镜像

执行以下命令，下载 docker 镜像。

```
docker pull ccr.ccs.tencentyun.com/qcloud/taco-train:ttf115-cu112-cvm-0.4.1
```

该镜像包含的软件版本信息如下：

- OS : 18.04.5
- python : 3.6.9
- cuda toolkits : V11.2.152
- cudnn library : 8.1.1
- nccl library : 2.8.4
- tencent-lightcc : 3.1.1
- HARP library : v1.3
- ttensorflow : 1.15.5

其中：

- LightCC 是云平台提供的基于 Horovod 深度定制优化的通信组件，完全兼容 Horovod API，不需要任何业务适配。
- HARP 是云平台提供的用户态协议栈，致力于提高 VPC 网络下的分布式训练的通信效率。以动态库的形式提供，官方 NCCL 初始化过程中会自动加载，不需要任何业务适配。
- ttensorflow 是云平台基于开源 tensorflow 1.15.5 添加了 CUDA 11 的支持，同时集成了 [TFRA](#)，用来支持动态 embedding 的特性。

## 启动 docker 镜像

执行以下命令，启动 docker 镜像。

```
docker run -it --rm --gpus all --privileged --net=host -v /sys:/sys -v /dev/hugepages:/dev/hugepages -v /usr/local/tfabric/tools:/usr/local/tfabric/tools ccr.ccs.tencentyun.com/qcloud/taco-train:ttf115-cu112-cvm-0.4.1
```

### 注意

/dev/hugepages 和 /usr/local/tfabric/tools 包含了 HARP 运行所需要的大页内存和配置文件。

## 分布式训练 benchmark 测试

## 说明

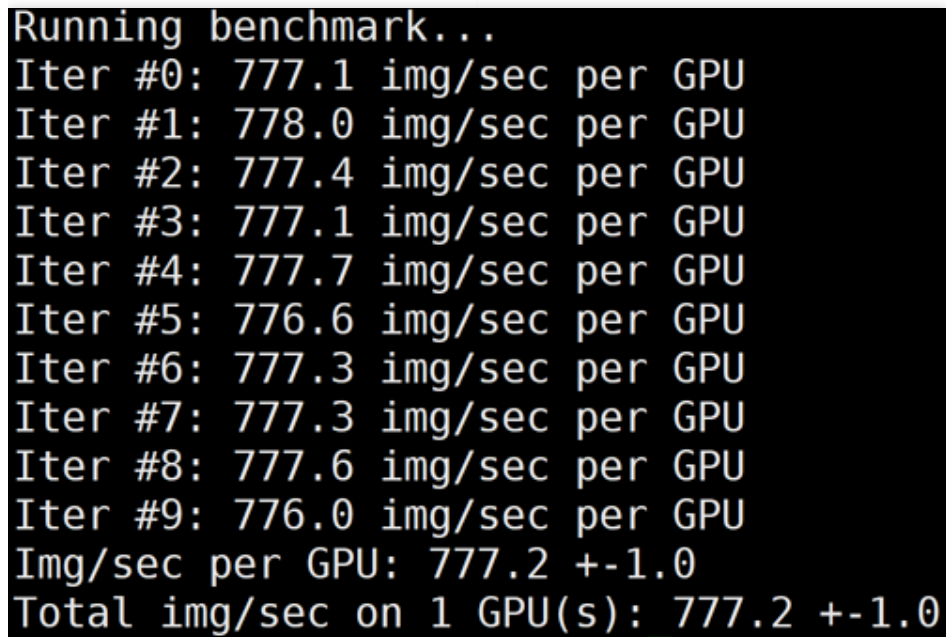
docker 镜像中的文件 `/mnt/tensorflow_synthetic_benchmark.py` 来自 [horovod example](#)。

## 单卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 1 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为A100的单卡 benchmark 结果：



```
Running benchmark...
Iter #0: 777.1 img/sec per GPU
Iter #1: 778.0 img/sec per GPU
Iter #2: 777.4 img/sec per GPU
Iter #3: 777.1 img/sec per GPU
Iter #4: 777.7 img/sec per GPU
Iter #5: 776.6 img/sec per GPU
Iter #6: 777.3 img/sec per GPU
Iter #7: 777.3 img/sec per GPU
Iter #8: 777.6 img/sec per GPU
Iter #9: 776.0 img/sec per GPU
Img/sec per GPU: 777.2 +/-1.0
Total img/sec on 1 GPU(s): 777.2 +/-1.0
```

## 单机多卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为 A100的单机8卡 benchmark 结果：

```
Running benchmark...
Iter #0: 761.9 img/sec per GPU
Iter #1: 764.0 img/sec per GPU
Iter #2: 763.2 img/sec per GPU
Iter #3: 763.8 img/sec per GPU
Iter #4: 763.2 img/sec per GPU
Iter #5: 763.0 img/sec per GPU
Iter #6: 763.8 img/sec per GPU
Iter #7: 763.0 img/sec per GPU
Iter #8: 762.9 img/sec per GPU
Iter #9: 762.8 img/sec per GPU
Img/sec per GPU: 763.1 +-1.1
Total img/sec on 8 GPU(s): 6105.2 +-9.0
```

### 多机多卡

1. 参考 开通实例 - 启动 docker 镜像 步骤，开通和配置多台训练机器。
2. 配置多台服务器 docker 间相互免密访问，详情请参见 [配置容器 SSH 免密访问](#)。
3. 执行以下命令，使用 TACO Train 进行多机训练加速。

```
/usr/local/openmpi/bin/mpirun -np 8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为 A100 的2机16卡 benchmark 结果：

```
Running benchmark...
Iter #0: 700.6 img/sec per GPU
Iter #1: 703.7 img/sec per GPU
Iter #2: 697.4 img/sec per GPU
Iter #3: 695.0 img/sec per GPU
Iter #4: 694.3 img/sec per GPU
Iter #5: 702.8 img/sec per GPU
Iter #6: 701.7 img/sec per GPU
Iter #7: 699.3 img/sec per GPU
Iter #8: 692.1 img/sec per GPU
Iter #9: 696.8 img/sec per GPU
Img/sec per GPU: 698.4 +-7.2
Total img/sec on 16 GPU(s): 11173.9 +-115.8
```

LightCC 的环境变量说明如下表：

环境变量	默认值	说明
LIGHT_2D_ALLREDUCE	0	是否使用2D-Allreduce 算法
LIGHT_INTRA_SIZE	8	2D-Allreduce 组内 GPU 数
LIGHT_HIERARCHICAL_THRESHOLD	1073741824	2D-Allreduce 的阈值，单位是字节，小于等于该阈值的数据才使用2D-Allreduce
LIGHT_TOPK_ALLREDUCE	0	是否使用 TOPK 压缩通信
LIGHT_TOPK_RATIO	0.01	使用 TOPK 压缩的比例
LIGHT_TOPK_THRESHOLD	1048576	TOPK 压缩的阈值，单位是字节，大于等于该阈值的数据才使用 TOPK 压缩通信
LIGHT_TOPK_FP16	0	压缩通信的 value 是否转成 FP16

4. 执行以下命令，关闭 TACO LightCC 加速进行测试。

```
# 修改环境变量，使用Horovod进行多机Allreduce
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by slot -x
NCCL_ALGO=RING -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python
3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为 A100的2机16卡，关闭 LightCC 之后的 benchmark 结果：

```
Running benchmark...
Iter #0: 494.7 img/sec per GPU
Iter #1: 496.1 img/sec per GPU
Iter #2: 492.7 img/sec per GPU
Iter #3: 490.2 img/sec per GPU
Iter #4: 488.7 img/sec per GPU
Iter #5: 487.8 img/sec per GPU
Iter #6: 494.6 img/sec per GPU
Iter #7: 487.3 img/sec per GPU
Iter #8: 488.9 img/sec per GPU
Iter #9: 490.1 img/sec per GPU
Img/sec per GPU: 491.1 +-5.9
Total img/sec on 16 GPU(s): 7857.7 +-94.1
```

5. 执行以下命令，同时关闭 LightCC 和 HARP 加速进行测试。

```
# 将HARP加速库rename为bak.libnccl-net.so即可关闭HARP加速。
/usr/local/openmpi/bin/mpirun -np 2 -H gpu1:1,gpu2:1 --allow-run-as-root -bind-to none -map-by slot mv
/usr/lib/x86_64-linux-gnu/libnccl-net.so /usr/lib/x86_64-linux-gnu/bak.libnccl-net.so

# 修改环境变量，使用Horovod进行多机Allreduce
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by slot -x
NCCL_ALGO=RING -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python
3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为 A100的2机16卡，同时关闭 LightCC 和 HARP 之后的 benchmark 结果：

```
Running benchmark...
Iter #0: 328.6 img/sec per GPU
Iter #1: 338.9 img/sec per GPU
Iter #2: 335.0 img/sec per GPU
Iter #3: 351.9 img/sec per GPU
Iter #4: 356.1 img/sec per GPU
Iter #5: 344.0 img/sec per GPU
Iter #6: 355.1 img/sec per GPU
Iter #7: 348.2 img/sec per GPU
Iter #8: 336.6 img/sec per GPU
Iter #9: 345.7 img/sec per GPU
Img/sec per GPU: 344.0 +-17.0
Total img/sec on 16 GPU(s): 5504.3 +-271.7
```

注意：

测试完如需恢复 HARP 加速能力，只需要把所有机器上的 bak.libnccl-net.so 重新命名为 libnccl-net.so 即可。

## 总结

本文测试数据如下：

机器：( A100 \* 8 ) + 50G VPC  
 容器：ccr.ccs.tencentyun.com/qcloud/taco-train:tff115-cu112-cvm-0.4.1  
 网络模型：ResNet50Batch：256  
 数据：synthetic data

GPU类型	#GPUs	Horovod+TCP		Horovod+HARP		LightCC+HARP	
		性能 ( img/sec )	线性加速比	性能 ( img/sec )	线性加速比	性能 ( img/sec )	线性加速比
A100	1	777	-	777	-	777	-

	8	6105	98.21%	6105	98.21%	6105	98.21%
	16	5504	44.27%	7857	63.20%	11173	89.87%

说明如下：

- 对于 A100 GPU机型，相比开源方案，使用 TACO 分布式训练加速组件之后，16卡A100的线性加速比从44.27%提升到89.87%，效果非常显著。
- LightCC 和 HARP 只在多机分布式训练当中才有加速效果，单机8卡场景由于 NVLink 的高速带宽存在，一般不需要额外的加速就能达到比较高的线性加速比。
- 上述 benchmark 脚本也支持除 ResNet50之外的其他模型，ModelName 请参考 [Keras Applications](#)。
- 上述 docker 镜像仅用于 demo，若您具备开发或者部署环境，请提供 OS/python/CUDA/tensorflow 版本信息，并联系我们售后提供特定版本的 TACO 加速组件。

# 在裸金属服务器上部署 TensorFlow 分布式训练集群

## 操作场景

本文介绍如何基于裸金属服务器搭建 Tensorflow+Taco Train 分布式训练集群。

## 操作步骤

### 开通实例

开通实例，其中实例、存储及镜像请参考以下信息选择，其余配置请参考 [通过开通页创建实例](#) 按需选择。

- 实例：选择 裸金属GPU实例。
- 系统盘：配置容量不小于50GB的云硬盘。您也可在创建实例后使用文件存储，详情参见 [在 Linux 客户端上使用 CFS 文件系统](#)。
- 镜像：建议选择公共镜像，公共镜像当中已安装 RDMA 网卡驱动。若您选择自定义镜像，则需要自行安装 RDMA 网卡驱动和 GPU 驱动。
- 操作系统 CentOS 7.6

### 配置实例环境

#### 验证 GPU 驱动

1. 登录实例。
2. 执行以下命令，验证 GPU 驱动是否安装成功。

```
nvdi-a-smi
```

查看输出结果是否为 GPU 状态：

- 是，代表 GPU 驱动安装成功。
- 否，请参考 [NVIDIA Driver Installation Quickstart Guide](#) 进行安装。

#### 安装 nv\_peer\_mem (可选)

多机通信的过程中，GPU 显存中的数据需要首先拷贝到内存中，然后通过网卡发出。通过 [GPU Direct RDMA 协](#)

议，可利用 GPU 和网卡直接通过 PCIe 进行 Peer2Peer 的数据交换这条更快速的路径，无需借助内存来进行数据的传递。

如需使用 GDR 进行数据传输，请在实例中执行以下命令，安装如下驱动。

```
git clone https://github.com/Mellanox/nv_peer_memory.git
cd ./nv_peer_memory/ && git checkout 1.0-9
make && insmod ./nv_peer_mem.ko
// 如果服务器发生了重启，nv_peer_mem驱动需要重新insmod
```

## 安装 docker 和 nvidia docker

1. 参见 [使用标准登录方式登录 Linux 实例](#)，登录实例。
2. 执行以下命令，安装 docker。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-docker.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 Docker 官方文档 [Install Docker Engine](#) 进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
+ sh -c 'yum install -y -q docker-ce-rootless-extras'
Package docker-ce-rootless-extras-20.10.16-3.el7.x86_64 already installed and latest version

=====

To run Docker as a non-privileged user, consider setting up the
Docker daemon in rootless mode for your user:

    dockerd-rootless-setuptool.sh install

Visit https://docs.docker.com/go/rootless/ to learn about rootless mode.

To run the Docker daemon as a fully privileged service, but granting non-root
users access, refer to https://docs.docker.com/go/daemon-access/

WARNING: Access to the remote API on a privileged Docker daemon is equivalent
to root access on the host. Refer to the 'Docker daemon attack surface'
documentation for details: https://docs.docker.com/go/attack-surface/

=====
```

3. 执行以下命令，安装 nvidia-docker2。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-nvidia-docker2.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 NVIDIA 官方文档 [Installation Guide & mdash](#) 进行安

装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
Running transaction
Installing : libnvidia-container1-1.9.0-1.x86_64
Installing : libnvidia-container-tools-1.9.0-1.x86_64
Installing : nvidia-container-toolkit-1.9.0-1.x86_64
Installing : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container-tools-1.9.0-1.x86_64
Verifying : nvidia-container-toolkit-1.9.0-1.x86_64
Verifying : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container1-1.9.0-1.x86_64

Installed:
nvidia-docker2.noarch 0:2.10.0-1

Dependency Installed:
libnvidia-container-tools.x86_64 0:1.9.0-1          libnvidia-container1.x86_64 0:1.9.0-1

Complete!
```

## 下载 docker 镜像

执行以下命令，下载 docker 镜像。

```
docker pull ccr.ccs.tencentyun.com/qcloud/taco-train:tff115-cu112-bm-0.4.2
```

该镜像包含的软件版本信息如下：

- OS : 18.04.5
- ofed: MLNX\_OFED\_LINUX-5.1-2.5.8.0
- python : 3.6.9
- cuda toolkits : V11.2.152
- cudnn library : 8.1.1
- nccl library : 2.8.4
- tencent-lightcc : 3.1.1
- ttensorflow : 1.15.5

其中：

- LightCC 是云平台提供的基于 Horovod 深度定制优化的通信组件，完全兼容 Horovod API，不需要任何业务适配。
- ttensorflow 是云平台基于开源 tensorflow 1.15.5 添加了 CUDA 11 的支持，同时集成了 [TFRA](#)，用来支持动态 embedding 的特性。

## 启动 docker 镜像

执行以下命令，启动 docker 镜像。

```
docker run -itd --rm --gpus all --shm-size=32g --ulimit memlock=-1 --ulimit stack=67108864 --net=host --privileged ccr.ccs.tencentyun.com/qcloud/taco-train:tff115-cu112-bm-0.4.2
```

### 注意

--privileged 选项使容器能够访问主机上的 RDMA 设备。

## 分布式训练 benchmark 测试

### 说明

docker 镜像中的文件 /mnt/tensorflow\_synthetic\_benchmark.py 来自 [horovod example](#)。

### 单卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 1 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_INDEX=3 -x NCCL_NET_GDR_LEVEL=0 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0 -mca btl^openib python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为A100的单卡 benchmark 结果：

```
Running benchmark...
Iter #0: 778.4 img/sec per GPU
Iter #1: 778.5 img/sec per GPU
Iter #2: 778.3 img/sec per GPU
Iter #3: 778.4 img/sec per GPU
Iter #4: 778.3 img/sec per GPU
Iter #5: 778.5 img/sec per GPU
Iter #6: 778.5 img/sec per GPU
Iter #7: 778.5 img/sec per GPU
Iter #8: 778.4 img/sec per GPU
Iter #9: 778.6 img/sec per GPU
Img/sec per GPU: 778.4 +-0.2
Total img/sec on 1 GPU(s): 778.4 +-0.2
```

### 单机多卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_INDEX=3 -x NCCL_N
```

```
ET_GDR_LEVEL=0 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0 -mca btl  
^openib python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=256
```

下图为 A100的8卡 benchmark 结果：

```
Running benchmark...  
Iter #0: 763.6 img/sec per GPU  
Iter #1: 761.9 img/sec per GPU  
Iter #2: 763.5 img/sec per GPU  
Iter #3: 763.3 img/sec per GPU  
Iter #4: 761.4 img/sec per GPU  
Iter #5: 764.0 img/sec per GPU  
Iter #6: 763.7 img/sec per GPU  
Iter #7: 763.5 img/sec per GPU  
Iter #8: 762.9 img/sec per GPU  
Iter #9: 762.7 img/sec per GPU  
Img/sec per GPU: 763.0 +-1.6  
Total img/sec on 8 GPU(s): 6104.4 +-12.6
```

### 多机多卡

1. 参考 [开通实例 - 启动 docker 镜像 步骤](#)，开通和配置多台训练机器。
2. 配置多台服务器 docker 间相互免密访问，详情请参见 [配置容器 SSH 免密访问](#)。
3. 执行以下命令，使用 TACO Train 进行多机训练加速。

```
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s  
lot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_I  
NDEX=3 -x NCCL_NET_GDR_LEVEL=0 -x HOROVOD_FUSION_THRESHOLD=0 -x HOROVOD_CYCLE_TI  
ME=0 -x LIGHT_INTRA_SIZE=8 -x LIGHT_2D_ALLREDUCE=1 -x LIGHT_TOPK_ALLREDUCE=1 -x LIGHT_T  
OPK_THRESHOLD=2097152 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0  
-mca btl ^openib python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=  
256
```

下图为 HCCPNV4h/A100 2机16卡 benchmark 结果：

```

Running benchmark...
Iter #0: 736.0 img/sec per GPU
Iter #1: 735.0 img/sec per GPU
Iter #2: 735.3 img/sec per GPU
Iter #3: 736.9 img/sec per GPU
Iter #4: 739.5 img/sec per GPU
Iter #5: 737.5 img/sec per GPU
Iter #6: 735.6 img/sec per GPU
Iter #7: 737.1 img/sec per GPU
Iter #8: 733.9 img/sec per GPU
Iter #9: 735.2 img/sec per GPU
Img/sec per GPU: 736.2 +-3.0
Total img/sec on 16 GPU(s): 11779.1 +-47.6

```

LightCC 的环境变量说明如下表：

环境变量	默认值	说明
LIGHT_2D_ALLREDUCE	0	是否使用2D-Allreduce 算法
LIGHT_INTRA_SIZE	8	2D-Allreduce 组内 GPU 数
LIGHT_HIERARCHICAL_THRESHOLD	1073741824	2D-Allreduce 的阈值，单位是字节，小于等于该阈值的数据才使用2D-Allreduce
LIGHT_TOPK_ALLREDUCE	0	是否使用 TOPK 压缩通信
LIGHT_TOPK_RATIO	0.01	使用 TOPK 压缩的比例
LIGHT_TOPK_THRESHOLD	1048576	TOPK 压缩的阈值，单位是字节，大于等于该阈值的数据才使用 TOPK 压缩通信
LIGHT_TOPK_FP16	0	压缩通信的 value 是否转成 FP16

4. 执行以下命令，关闭 TACO LightCC 加速进行测试。

```
# 去掉LIGHT_xx的环境变量，即可使用Horovod进行多机Allreduce
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_I
NDEX=3 -x NCCL_NET_GDR_LEVEL=0 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_inclu
de bond0 -mca btl ^openib python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --b
atch-size=256
```

下图为 A100 2机16卡，关闭 LightCC 之后的 benchmark 结果：

```
Running benchmark...
Iter #0: 673.6 img/sec per GPU
Iter #1: 678.8 img/sec per GPU
Iter #2: 669.6 img/sec per GPU
Iter #3: 677.4 img/sec per GPU
Iter #4: 671.3 img/sec per GPU
Iter #5: 672.9 img/sec per GPU
Iter #6: 676.1 img/sec per GPU
Iter #7: 680.1 img/sec per GPU
Iter #8: 674.2 img/sec per GPU
Iter #9: 670.0 img/sec per GPU
Img/sec per GPU: 674.4 +-6.8
Total img/sec on 16 GPU(s): 10790.4 +-108.2
```

### 多机多卡GDR

执行以下命令，使用 GDR 进行测试。

注意：

使用 GDR 需安装 `nv_peer_mem`，详情请参见 [安装 nv\\_peer\\_mem](#)。

```
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_I
NDEX=3 -x NCCL_NET_GDR_LEVEL=2 -x HOROVOD_FUSION_THRESHOLD=0 -x HOROVOD_CYCLE_TI
ME=0 -x LIGHT_INTRA_SIZE=8 -x LIGHT_2D_ALLREDUCE=1 -x LIGHT_TOPK_ALLREDUCE=1 -x LIGHT_T
OPK_THRESHOLD=2097152 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0
-mca btl ^openib python3 /mnt/tensorflow_synthetic_benchmark.py --model=ResNet50 --batch-size=
256
```

GDR 通常在大模型或者集群规模较大时有显著的加速效果，测试结果如下图所示：

```
Running benchmark...
Iter #0: 733.5 img/sec per GPU
Iter #1: 731.3 img/sec per GPU
Iter #2: 735.2 img/sec per GPU
Iter #3: 734.1 img/sec per GPU
Iter #4: 733.7 img/sec per GPU
Iter #5: 734.8 img/sec per GPU
Iter #6: 734.5 img/sec per GPU
Iter #7: 736.3 img/sec per GPU
Iter #8: 733.8 img/sec per GPU
Iter #9: 733.7 img/sec per GPU
Img/sec per GPU: 734.1 +-2.4
Total img/sec on 16 GPU(s): 11745.3 +-38.6
```

## 总结

本文使用环境及测试数据如下：

机器：( A100 \* 8 ) + 100G RDMA + 25G VPC  
 容器：ccr.ccs.tencentyun.com/qcloud/taco-train:ttf115-cu112-bm-0.4.2  
 网络模型：ResNet50Batch：256  
 数据：synthetic data

GPU卡类型	#GPUs	Horovod+RDMA		LightCC+RDMA	
		性能 ( img/sec )	线性加速比	性能 ( img/sec )	线性加速比
A100	1	778	-	778	-
	8	6104	98.07%	6104	98.07%
	16	10790	86.68%	11779	94.63%

说明如下：

- 对于 A100裸金属机型，相比开源方案，2机16卡通过 LightCC 可以将线性加速比从86.68%提升到94.63%。
- 上述 benchmark 也支持除 ResNet50之外的其他模型，ModelName 请参考 [Keras Applications](#)。
- 上述 docker 镜像仅用于 demo，若您需使用自己的 docker 开发环境，请参考 [容器安装用户态 RDMA 驱动](#) 安装网卡驱动。
- 如需特定 OS/python/CUDA/tensorflow 版本的 LightCC 加速组件，请联系工程师。



# 在 CVM 上部署 PyTorch 分布式训练集群

## 操作场景

本文介绍如何基于云服务器 CVM 搭建 torch+Taco Train 分布式训练集群。

## 操作步骤

### 开通实例

开通实例，其中实例、存储及镜像请参考以下信息选择，其余配置请参考 [通过开通页创建实例](#) 按需选择。

- 实例：选择合适的计算实例。
- 系统盘：配置容量不小于50GB的云硬盘。您也可在创建实例后使用文件存储，详情参见 [在 Linux 客户端上使用 CFS 文件系统](#)。
- 镜像：建议选择公共镜像，您也可选择自定义镜像。
- 操作系统请使用 CentOS 8.0/CentOS 7.8/Ubuntu 20.04/Ubuntu 18.04。

### 配置实例环境

#### 验证 GPU 驱动

1. 参考 [使用标准登录方式登录 Linux 实例](#)，登录实例。
2. 执行以下命令，验证 GPU 驱动是否安装成功。

```
nvdi-a-smi
```

查看输出结果是否为 GPU 状态：

- 是，代表 GPU 驱动安装成功。
- 否，请参考 [NVIDIA Driver Installation Quickstart Guide](#) 进行安装。

### 配置 HARP 分布式训练环境

1. 参考 [配置 HARP 分布式训练环境](#)，配置所需环境。
2. 配置完成后，执行以下命令进行验证，若配置文件存在，则表示已配置成功。

```
ls /usr/local/tfabric/tools/config/ztcp*.conf
```

## 安装 docker 和 nvidia docker

1. 执行以下命令，安装 docker。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-docker.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 Docker 官方文档 [Install Docker Engine](#) 进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
+ sh -c 'yum install -y -q docker-ce-rootless-extras'
Package docker-ce-rootless-extras-20.10.16-3.el7.x86_64 already installed and latest version
=====

To run Docker as a non-privileged user, consider setting up the
Docker daemon in rootless mode for your user:

    dockerd-rootless-setuptool.sh install

Visit https://docs.docker.com/go/rootless/ to learn about rootless mode.

To run the Docker daemon as a fully privileged service, but granting non-root
users access, refer to https://docs.docker.com/go/daemon-access/

WARNING: Access to the remote API on a privileged Docker daemon is equivalent
to root access on the host. Refer to the 'Docker daemon attack surface'
documentation for details: https://docs.docker.com/go/attack-surface/
=====
```

2. 执行以下命令，安装 nvidia-docker2。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-nvidia-docker2.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 NVIDIA 官方文档 [Installation Guide & mdash](#) 进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
Running transaction
Installing : libnvidia-container1-1.9.0-1.x86_64
Installing : libnvidia-container-tools-1.9.0-1.x86_64
Installing : nvidia-container-toolkit-1.9.0-1.x86_64
Installing : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container-tools-1.9.0-1.x86_64
Verifying : nvidia-container-toolkit-1.9.0-1.x86_64
Verifying : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container1-1.9.0-1.x86_64

Installed:
nvidia-docker2.noarch 0:2.10.0-1

Dependency Installed:
libnvidia-container-tools.x86_64 0:1.9.0-1      libnvidia-container1.x86_64 0:1.9.0-1      nvidia-container-toolkit.x86_64 0:1.9.0-1

Complete!
```

## 下载 docker 镜像

执行以下命令，下载 docker 镜像。

```
docker pull ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-cvm-0.4.3
```

该镜像包含的软件版本信息如下：

- OS : Ubuntu 20.04.4 LTS
- python : 3.8.10
- cuda toolkits : V11.3.109
- cudnn library : 8.2.0
- nccl library : 2.9.9
- tencent-lightcc : 3.1.1
- HARP library : v1.3
- torch : 1.11.0+cu113

其中：

- LightCC 是云平台提供的基于 Horovod 深度定制优化的通信组件，完全兼容 Horovod API，不需要任何业务适配。
- HARP 是云平台提供的用户态协议栈，致力于提高 VPC 网络下的分布式训练的通信效率。以动态库的形式提供，官方 NCCL 初始化过程中会自动加载，不需要任何业务适配。
- torch 是官方版本1.11.0版本。

## 启动 docker 镜像

执行以下命令，启动 docker 镜像。

```
docker run -it --rm --gpus all --privileged --net=host -v /sys:/sys -v /dev/hugepages:/dev/hugepages -v /usr/local/tfabric/tools:/usr/local/tfabric/tools ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-cvm-0.4.3
```

### 注意

`/dev/hugepages` 和 `/usr/local/tfabric/tools` 包含了 HARP 运行所需要的大页内存和配置文件。

## 分布式训练 benchmark 测试

### 说明

docker 镜像中的文件 `/mnt/tensorflow_synthetic_benchmark.py` 来自 [horovod example](#)。

### 单卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 1 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_benchmark.py --model=vgg16 --batch-size=128
```

下图为 V100的单卡 benchmark 结果：

```
Model: vgg16
Batch size: 128
Number of GPUs: 1
Running warmup...
Running benchmark...
Iter #0: 253.0 img/sec per GPU
Iter #1: 252.1 img/sec per GPU
Iter #2: 252.7 img/sec per GPU
Iter #3: 251.2 img/sec per GPU
Iter #4: 251.7 img/sec per GPU
Iter #5: 251.1 img/sec per GPU
Iter #6: 251.1 img/sec per GPU
Iter #7: 250.9 img/sec per GPU
Iter #8: 250.3 img/sec per GPU
Iter #9: 250.8 img/sec per GPU
Img/sec per GPU: 251.5 +-1.6
Total img/sec on 1 GPU(s): 251.5 +-1.6
```

### 单机多卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_benchmark.py --model=vgg16 --batch-size=128
```

下图为 V100的单机8卡 benchmark 结果：

```
Model: vgg16
Batch size: 128
Number of GPUs: 8
Running warmup...
Running benchmark...
Iter #0: 248.0 img/sec per GPU
Iter #1: 248.0 img/sec per GPU
Iter #2: 247.3 img/sec per GPU
Iter #3: 247.3 img/sec per GPU
Iter #4: 247.5 img/sec per GPU
Iter #5: 247.0 img/sec per GPU
Iter #6: 246.9 img/sec per GPU
Iter #7: 246.7 img/sec per GPU
Iter #8: 246.9 img/sec per GPU
Iter #9: 246.3 img/sec per GPU
Img/sec per GPU: 247.2 +-1.0
Total img/sec on 8 GPU(s): 1977.4 +-8.2
```

## 多机多卡

1. 参考 [开通实例 - 启动 docker 镜像 步骤](#)，开通和配置多台训练机器。
2. 配置多台服务器 docker 间相互免密访问，详情请参见 [配置容器 SSH 免密访问](#)。
3. 执行以下命令，使用 TACO Train 进行多机训练加速。

```
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_ALGO=RING -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x HOROVOD_MPI_T
HREADS_DISABLE=1 -x HOROVOD_FUSION_THRESHOLD=0 -x HOROVOD_CYCLE_TIME=0 -x LIGHT_I
NTRA_SIZE=8 -x LIGHT_2D_ALLREDUCE=1 -x LIGHT_TOPK_ALLREDUCE=1 -x LIGHT_TOPK_THRESHOLD
=2097152 -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_
benchmark.py --model=vgg16 --batch-size=128
```

下图为 V100的2机16卡 benchmark 结果：

```
Running benchmark...
Iter #0: 222.9 img/sec per GPU
Iter #1: 225.2 img/sec per GPU
Iter #2: 222.9 img/sec per GPU
Iter #3: 225.3 img/sec per GPU
Iter #4: 221.7 img/sec per GPU
Iter #5: 219.5 img/sec per GPU
Iter #6: 224.5 img/sec per GPU
Iter #7: 221.6 img/sec per GPU
Iter #8: 221.7 img/sec per GPU
Iter #9: 221.1 img/sec per GPU
Img/sec per GPU: 222.6 +-3.5
Total img/sec on 16 GPU(s): 3562.2 +-55.9
```

LightCC 的环境变量说明如下表：

环境变量	默认值	说明
LIGHT_2D_ALLREDUCE	0	是否使用2D-Allreduce 算法
LIGHT_INTRA_SIZE	8	2D-Allreduce 组内 GPU 数
LIGHT_HIERARCHICAL_THRESHOLD	1073741824	2D-Allreduce 的阈值，单位是字节，小于等于该阈值的数据才使用2D-Allreduce
LIGHT_TOPK_ALLREDUCE	0	是否使用 TOPK 压缩通信
LIGHT_TOPK_RATIO	0.01	使用 TOPK 压缩的比例
LIGHT_TOPK_THRESHOLD	1048576	TOPK 压缩的阈值，单位是字节，大于等于该阈值的数据才使用 TOPK 压缩通信
LIGHT_TOPK_FP16	0	压缩通信的 value 是否转成 FP16

4. 执行以下命令，关闭 TACO LightCC 加速进行测试。

```
# 修改环境变量，使用Horovod进行多机Allreduce
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_ALGO=RING -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x LD_LIBRARY_PATH
-x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_benchmark.py --model=vgg16
--batch-size=128
```

下图为 V100的2机16卡，关闭 LightCC 后的 benchmark 结果：

```
Model: vgg16
Batch size: 128
Number of GPUs: 16
Running warmup...
Running benchmark...
Iter #0: 103.4 img/sec per GPU
Iter #1: 103.4 img/sec per GPU
Iter #2: 103.7 img/sec per GPU
Iter #3: 103.2 img/sec per GPU
Iter #4: 103.8 img/sec per GPU
Iter #5: 103.6 img/sec per GPU
Iter #6: 104.0 img/sec per GPU
Iter #7: 103.2 img/sec per GPU
Iter #8: 103.4 img/sec per GPU
Iter #9: 95.8 img/sec per GPU
Img/sec per GPU: 102.8 +-4.6
Total img/sec on 16 GPU(s): 1644.1 +-73.5
```

5. 执行以下命令，同时关闭 LightCC 和 HARP 加速进行测试。

```
# 将HARP加速库rename为bak.libnccl-net.so即可关闭HARP加速。/usr/local/openmpi/bin/mpirun -np 2 -H gpu1:1,gpu2:1 --allow-run-as-root -bind-to none -map-by slot mv /usr/lib/x86_64-linux-gnu/libnccl-net.so /usr/lib/x86_64-linux-gnu/bak.libnccl-net.so# 修改环境变量，使用Horovod进行多机Allreduce /usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_ALGO=RING -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_benchmark.py --model=vgg16 --batch-size=128
```

下图为 V100的2机16卡，同时关闭 LightCC 和 HARP 后的 benchmark 结果：

```

Model: vgg16
Batch size: 128
Number of GPUs: 16
Running warmup...
Running benchmark...
Iter #0: 52.8 img/sec per GPU
Iter #1: 51.4 img/sec per GPU
Iter #2: 51.6 img/sec per GPU
Iter #3: 53.7 img/sec per GPU
Iter #4: 54.1 img/sec per GPU
Iter #5: 54.0 img/sec per GPU
Iter #6: 53.0 img/sec per GPU
Iter #7: 54.1 img/sec per GPU
Iter #8: 52.9 img/sec per GPU
Iter #9: 52.2 img/sec per GPU
Img/sec per GPU: 53.0 +-1.9
Total img/sec on 16 GPU(s): 847.8 +-30.1

```

注意：

测试完如需恢复 HARP 加速能力，只需要把所有机器上的 bak.libnccl-net.so 重新命名为 libncc-net.so 即可。

## 总结

本文测试数据如下：

机器：( V100 \* 8 ) + 25G VPC  
 容器：ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-cvm-0.4.1  
 网络模型：VGG16Batch：128  
 数据：synthetic data

GPU类型	#GPUs	Horovod+TCP		Horovod+HARP		LightCC+HARP	
		性能 ( img/ sec )	线性加速比	性能 ( img/ sec )	线性加速比	性能 ( img/ sec )	线性加速比
V100	1	251	-	251	-	251	-
	8	1977	98%	1977	98%	1977	98%
	16	847	21%	1644	40%	3562	88%

说明如下：

- 对于 V100机型，相比开源方案，使用 TACO 分布式训练加速组件之后，16卡V100的线性加速比从93%提升到 97%（由于当前网络模型的加速已经很高，软件优化提升的空间有限）。
- LightCC 和 HARP 只在多机分布式训练当中才有加速效果，单机8卡场景由于 NVLink 的高速带宽存在，一般不需要额外的加速就能达到比较高的线性加速比。
- 上述 benchmark 脚本也可以支持除了 VGG16之外的其他模型。ModelName 请参考 [Keras Applications](#)。
- 上述 docker 镜像仅用于 demo，若您具备开发或者部署环境，请提供 OS/python/CUDA/torch 版本信息，来获取特定版本的 TACO 加速组件。

# 在裸金属服务器上部署 PyTorch 分布式训练集群

## 操作场景

本文介绍如何基于裸金属服务器搭建 torch+Taco Train 分布式训练集群。

## 操作步骤

### 开通实例

开通实例，其中实例、存储及镜像请参考以下信息选择，其余配置请参考 [通过开通页创建实例](#) 按需选择。

- 实例：选择 裸金属GPU实例。
- 系统盘：配置容量不小于50GB的云硬盘。您也可在创建实例后使用文件存储，详情参见 [在 Linux 客户端上使用 CFS 文件系统](#)。
- 镜像：建议选择公共镜像，公共镜像当中已安装 RDMA 网卡驱动。若您选择自定义镜像，则需要自行安装 RDMA 网卡驱动和 GPU 驱动。
- 操作系统 CentOS 7.6

### 配置实例环境

#### 验证 GPU 驱动

1. 登录实例。
2. 执行以下命令，验证 GPU 驱动是否安装成功。

```
nvdiia-smi
```

查看输出结果是否为 GPU 状态：

- 是，代表 GPU 驱动安装成功。
- 否，请参考 [NVIDIA Driver Installation Quickstart Guide](#) 进行安装。

#### 安装 nv\_peer\_mem (可选)

多机通信的过程中，GPU 显存中的数据需要首先拷贝到内存中，然后通过网卡发出。通过 [GPU Direct RDMA 协](#)

议，可利用 GPU 和网卡直接通过 PCIe 进行 Peer2Peer 的数据交换这条更快速的路径，无需借助内存来进行数据的传递。

如需使用 GDR 进行数据传输，请在实例中执行以下命令，安装如下驱动。

```
git clone https://github.com/Mellanox/nv_peer_memory.git
cd ./nv_peer_memory/ && git checkout 1.0-9
make && insmod ./nv_peer_mem.ko
// 如果服务器发生了重启，nv_peer_mem驱动需要重新insmod
```

## 安装 docker 和 nvidia docker

1. 参考 [使用标准登录方式登录 Linux 实例](#)，登录实例。

2. 执行以下命令，安装 docker。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-docker.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 Docker 官方文档 [Install Docker Engine](#) 进行安装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
+ sh -c 'yum install -y -q docker-ce-rootless-extras'
Package docker-ce-rootless-extras-20.10.16-3.el7.x86_64 already installed and latest version

=====

To run Docker as a non-privileged user, consider setting up the
Docker daemon in rootless mode for your user:

    dockerd-rootless-setuptool.sh install

Visit https://docs.docker.com/go/rootless/ to learn about rootless mode.

To run the Docker daemon as a fully privileged service, but granting non-root
users access, refer to https://docs.docker.com/go/daemon-access/

WARNING: Access to the remote API on a privileged Docker daemon is equivalent
to root access on the host. Refer to the 'Docker daemon attack surface'
documentation for details: https://docs.docker.com/go/attack-surface/

=====
```

3. 执行以下命令，安装 nvidia-docker2。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/get-nvidia-docker2.sh | sudo bash
```

若您无法通过该命令安装，请尝试多次执行命令，或参考 NVIDIA 官方文档 [Installation Guide & mdash](#) 进行安

装。

本文以 CentOS 为例，安装成功后，返回结果如下图所示：

```
Running transaction
Installing : libnvidia-container1-1.9.0-1.x86_64
Installing : libnvidia-container-tools-1.9.0-1.x86_64
Installing : nvidia-container-toolkit-1.9.0-1.x86_64
Installing : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container-tools-1.9.0-1.x86_64
Verifying : nvidia-container-toolkit-1.9.0-1.x86_64
Verifying : nvidia-docker2-2.10.0-1.noarch
Verifying : libnvidia-container1-1.9.0-1.x86_64

Installed:
nvidia-docker2.noarch 0:2.10.0-1

Dependency Installed:
libnvidia-container-tools.x86_64 0:1.9.0-1      libnvidia-container1.x86_64 0:1.9.0-1      nvidia-container-toolkit.x86_64 0:1.9.0-1

Complete!
```

## 下载 docker 镜像

执行以下命令，下载 docker 镜像。

```
docker pull ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-bm-0.4.1
```

该镜像包含的软件版本信息如下：

- OS : Ubuntu 20.04.4 LTS
- ofed: MLNX\_OFED\_LINUX-5.4-3.1.0.0
- python : 3.8.10
- cuda toolkits : V11.3.109
- cudnn library : 8.2.0
- nccl library : 2.9.9
- tencent-lightcc : 3.1.1
- torch : 1.11.0+cu113

其中：

- LightCC 是云平台提供的基于 Horovod 深度定制优化的通信组件，完全兼容 Horovod API，不需要任何业务适配。
- torch 为官方版本。

## 启动 docker 镜像

执行以下命令，启动 docker 镜像。

```
docker run -itd --rm --gpus all --shm-size=32g --ulimit memlock=-1 --ulimit stack=67108864 --net=host --privileged ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-bm-0.4.1
```

### 注意

--privileged 选项使容器能够访问主机上的 RDMA 设备。

## 分布式训练 benchmark 测试

### 说明

docker 镜像中的文件 /mnt/pytorch\_synthetic\_benchmark.py 来自 [horovod example](#)。

### 单卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 1 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_SOCKET_IFNAME=eth0 -x LD_LIBRARY_PATH -x PATH -mca btl_tcp_if_include eth0 python3 /mnt/pytorch_synthetic_benchmark.py --model=vgg16 --batch-size=128
```

下图为 V100的单卡 benchmark 结果：

```
Model: resnet50
Batch size: 256
Number of GPUs: 1
Running warmup...
Running benchmark...
Iter #0: 819.7 img/sec per GPU
Iter #1: 819.6 img/sec per GPU
Iter #2: 819.6 img/sec per GPU
Iter #3: 819.6 img/sec per GPU
Iter #4: 819.6 img/sec per GPU
Iter #5: 819.6 img/sec per GPU
Iter #6: 819.6 img/sec per GPU
Iter #7: 819.6 img/sec per GPU
Iter #8: 819.5 img/sec per GPU
Iter #9: 819.7 img/sec per GPU
Img/sec per GPU: 819.6 +-0.1
Total img/sec on 1 GPU(s): 819.6 +-0.1
```

### 单机多卡

执行以下命令，进行测试。

```
/usr/local/openmpi/bin/mpirun -np 8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_INDEX=3 -x NCCL_NET_GDR_LEVEL=0 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0 -mca btl^openib python3 /mnt/pytorch_synthetic_benchmark.py --model resnet50 --batch-size=256
```

下图为 A100的8卡 benchmark 结果：

```
Model: resnet50
Batch size: 256
Number of GPUs: 8
Running warmup...
Running benchmark...
Iter #0: 808.6 img/sec per GPU
Iter #1: 808.8 img/sec per GPU
Iter #2: 808.9 img/sec per GPU
Iter #3: 809.1 img/sec per GPU
Iter #4: 808.7 img/sec per GPU
Iter #5: 808.4 img/sec per GPU
Iter #6: 809.2 img/sec per GPU
Iter #7: 808.0 img/sec per GPU
Iter #8: 809.1 img/sec per GPU
Iter #9: 808.5 img/sec per GPU
Img/sec per GPU: 808.7 +-0.7
Total img/sec on 8 GPU(s): 6469.8 +-5.2
```

## 多机多卡

1. 参考 [开通实例 - 启动 docker 镜像 步骤](#)，开通和配置多台训练机器。
2. 配置多台服务器 docker 间相互免密访问，详情请参见 [配置容器 SSH 免密访问](#)。
3. 执行以下命令，使用 TACO Train 进行多机训练加速。

```
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_I
NDEX=3 -x NCCL_NET_GDR_LEVEL=0 -x HOROVOD_FUSION_THRESHOLD=0 -x HOROVOD_CYCLE_TI
ME=0 -x LIGHT_INTRA_SIZE=8 -x LIGHT_2D_ALLREDUCE=1 -x LIGHT_TOPK_ALLREDUCE=1 -x LIGHT_T
OPK_THRESHOLD=2097152 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0
-mca btl ^openib python3 /mnt/pytorch_synthetic_benchmark.py --model resnet50 --batch-size=256
```

下图为 A100 2机16卡 benchmark 结果：

```
Running benchmark...
Iter #0: 781.5 img/sec per GPU
Iter #1: 784.1 img/sec per GPU
Iter #2: 782.9 img/sec per GPU
Iter #3: 782.8 img/sec per GPU
Iter #4: 783.4 img/sec per GPU
Iter #5: 784.3 img/sec per GPU
Iter #6: 783.2 img/sec per GPU
Iter #7: 782.9 img/sec per GPU
Iter #8: 783.9 img/sec per GPU
Iter #9: 783.6 img/sec per GPU
Img/sec per GPU: 783.3 +-1.5
Total img/sec on 16 GPU(s): 12532.5 +-24.4
```

LightCC 的环境变量说明如下表：

环境变量	默认值	说明
LIGHT_2D_ALLREDUCE	0	是否使用2D-Allreduce 算法
LIGHT_INTRA_SIZE	8	2D-Allreduce 组内 GPU 数
LIGHT_HIERARCHICAL_THRESHOLD	1073741824	2D-Allreduce 的阈值，单位是字节，小于等于该阈值的数据才使用2D-Allreduce
LIGHT_TOPK_ALLREDUCE	0	是否使用 TOPK 压缩通信
LIGHT_TOPK_RATIO	0.01	使用 TOPK 压缩的比例
LIGHT_TOPK_THRESHOLD	1048576	TOPK 压缩的阈值，单位是字节，大于等于该阈值的数据才使用 TOPK 压缩通信
LIGHT_TOPK_FP16	0	压缩通信的 value 是否转成 FP16

4. 执行以下命令，关闭 TACO LightCC 加速进行测试。

```
# 去掉LIGHT_xx的环境变量，即可使用Horovod进行多机Allreduce/usr/local/openmpi/bin/mpirun -np 16
-H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by slot -x NCCL_DEBUG=INFO -x NCCL_IB_
DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_INDEX=3 -x NCCL_NET_GDR_LEVEL=0 -
x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0 -mca btl ^openib python3 /
mnt/pytorch_synthetic_benchmark.py --model resnet50 --batch-size=256
```

下图为 A100 2机16卡，关闭 LightCC 之后的 benchmark 结果：

```
Model: resnet50
Batch size: 256
Number of GPUs: 16
Running warmup...
Running benchmark...
Iter #0: 772.3 img/sec per GPU
Iter #1: 766.2 img/sec per GPU
Iter #2: 769.4 img/sec per GPU
Iter #3: 766.5 img/sec per GPU
Iter #4: 767.0 img/sec per GPU
Iter #5: 767.3 img/sec per GPU
Iter #6: 768.9 img/sec per GPU
Iter #7: 772.0 img/sec per GPU
Iter #8: 770.6 img/sec per GPU
Iter #9: 767.2 img/sec per GPU
Img/sec per GPU: 768.7 +-4.2
Total img/sec on 16 GPU(s): 12299.7 +-67.6
```

### 多机多卡GDR

执行以下命令，使用 GDR 进行测试。

注意：

使用 GDR 需安装 `nv_peer_mem`，详情请参见 [安装 nv\\_peer\\_mem](#)。

```
/usr/local/openmpi/bin/mpirun -np 16 -H gpu1:8,gpu2:8 --allow-run-as-root -bind-to none -map-by s
lot -x NCCL_DEBUG=INFO -x NCCL_IB_DISABLE=0 -x NCCL_SOCKET_IFNAME=bond0 -x NCCL_IB_GID_I
NDEX=3 -x NCCL_NET_GDR_LEVEL=2 -x HOROVOD_FUSION_THRESHOLD=0 -x HOROVOD_CYCLE_TI
ME=0 -x LIGHT_INTRA_SIZE=8 -x LIGHT_2D_ALLREDUCE=1 -x LIGHT_TOPK_ALLREDUCE=1 -x LIGHT_T
OPK_THRESHOLD=2097152 -x LD_LIBRARY_PATH -x PATH -mca pml ob1 -mca btl_tcp_if_include bond0
-mca btl ^openib python3 /mnt/pytorch_synthetic_benchmark.py --model resnet50 --batch-size=256
```

GDR 通常在大模型或者集群规模较大时有显著的加速效果，测试结果如下图所示：

```
Running benchmark...
Iter #0: 782.0 img/sec per GPU
Iter #1: 781.1 img/sec per GPU
Iter #2: 782.0 img/sec per GPU
Iter #3: 783.1 img/sec per GPU
Iter #4: 782.1 img/sec per GPU
Iter #5: 780.7 img/sec per GPU
Iter #6: 781.6 img/sec per GPU
Iter #7: 780.4 img/sec per GPU
Iter #8: 783.0 img/sec per GPU
Iter #9: 784.3 img/sec per GPU
Img/sec per GPU: 782.0 +-2.2
Total img/sec on 16 GPU(s): 12512.6 +-35.7
```

## 总结

本文使用环境及测试数据如下：

机器：( A100 \* 8 ) + 100G RDMA + 25G VPC  
 容器：ccr.ccs.tencentyun.com/qcloud/taco-train:torch111-cu113-bm-0.4.1  
 网络模型：ResNet50Batch：256  
 数据：synthetic data

GPU类型	#GPUs	Horovod+RDMA		LightCC+RDMA	
		性能 ( img/sec )	线性加速比	性能 ( img/sec )	线性加速比
A100	1	819	-	819	-
	8	6469	98.73%	6469	98.73%
	16	12299	93.85%	12532	95.63%

说明如下：

- 对于 A100，相比开源方案，2机16卡通过 LightCC 可将线性加速比从93.85%提升到95.63%。
- 上述 benchmark 也支持除 ResNet50之外的其他模型，ModelName 请参考 [Keras Applications](#)。
- 上述 docker 镜像仅用于 demo，若您需使用自己的 docker 开发环境，请参考 [容器安装用户态 RDMA 驱动](#) 安装网卡驱动。
- 如需特定 OS/python/CUDA/tensorflow 版本的 LightCC 加速组件，请联系工程师获取。

# 组件配置和使用

## TCCL 使用说明

### 操作场景

本文介绍如何在云平台环境中配置 TCCL 加速通信库，实现您在云平台 RDMA 环境中多机多卡通通信的性能提升。在大模型训练场景，对比开源的 NCCL 方案，TCCL 预计约可以提升 50% 带宽利用率。

### 操作步骤

#### 准备环境

- 1、创建 分别支持 1.6Tbps 和 800Gbps RDMA 网络的高性能计算集群实例。
- 2、为 GPU 型实例安装 GPU 驱动和 nvidia-fabricmanager 服务。

注意：

TCCL 运行软件环境要求 glibc 版本 2.17 以上，CUDA 版本 10.0 以上。

#### 选择安装方式

TCCL目前支持三种使用方式安装，您可以根据需要选择适合业务场景的安装方式使用。

- TCCL通信库 + 编译安装pytorch
- TCCL通信库 + pytorch通信插件
- NCCL插件 + 排序的IP列表

说明：

由于当前大模型训练基本都基于 Pytorch 框架，所以主要以 Pytorch 为例进行说明，

TCCL的三种接入方案对比如下表：

安装方式	方法一：编译安装 Pytorch	方法二：安装 Pytorch 通信插件	(推荐)方法三：安装 NCCL 通信插件
使用步骤	<ul style="list-style-type: none"><li>• 安装 TCCL</li><li>• 重新编译安装 Pytorch</li></ul>	<ul style="list-style-type: none"><li>• 安装 Pytorch 通信插件</li><li>• 修改分布式通信后端</li></ul>	<ul style="list-style-type: none"><li>• 安装 NCCL 插件</li><li>• 修改启动脚本</li></ul>
优点	对业务代码无入侵	安装方便	安装方便

缺点	需要重新编译安装 Pytorch 对软件环境有要求	需要修改业务代码 对软件环境有要求	集群节点扩充之后，需要更新排序列表
软件环境依赖	对应 NCCL 版本 2.12 要求 glibc 版本 2.17 以上 要求 CUDA 版本 10.0 以上	当前安装包仅支持 Pytorch 1.12 要求 glibc 版本 2.17 以上  要求 CUDA 版本 10.0 以上	安装 NCCL 即可

如果您的机器资源和模型训练场景相对比较固定，推荐使用方法3，兼容不同的NCCL版本和CUDA版本，安装使用方便，不需要修改业务代码或者重新编译pytorch。

如果您的资源需要提供给不同的业务团队，或者经常有扩容的需求，推荐使用前两种方法，不需要算法人员或者调度框架刻意去感知机器的网络拓扑信息。

如果您不希望对业务代码做适配，那么可以使用方法1，只需要重新编译pytorch框架。

## 配置 TCCL 环境并验证

### 方法一：编译安装 Pytorch

由于社区pytorch默认采用静态方式连接NCCL通信库，所以无法通过替换共享库的方式使用TCCL。

#### 1、安装TCCL

以 Ubuntu 20.04 为例，您可以使用以下命令安装，安装之后TCCL位于 /opt/tencent/tccl 目录。

```
# 卸载已有tccl版本和nccl插件
dpkg -r tccl && dpkg -r nccl-rdma-sharp-plugins

# 下载安装tccl v1.5版本
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/tccl/TCCL_1.5-ubuntu.20.04.5_amd64.deb && dpkg -i TCCL_1.5-ubuntu.20.04.5_amd64.deb
```

如果您使用 CentOS 或 TencentOS，参考以下步骤安装：

```
# 卸载已有tccl版本和nccl插件
rpm -e tccl && rpm -e nccl-rdma-sharp-plugins-1.0-1.x86_64

# 下载tccl v1.5版本
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/tccl/tccl-1.5-1.tl2.x86_64.rpm && rpm -ivh --nodeps --force tccl-1.5-1.tl2.x86_64.rpm && rm -f tccl-1.5-1.tl2.x86_64.rpm
```

#### 2、重新编译安装 Pytorch

以下为 Pytorch 源码安装示例，详情参考[官网 Pytorch 安装说明](#)。

```
#!/bin/bash

# 卸载当前版本
pip uninstall -y torch

# 下载pytorch源码
git clone --recursive https://github.com/pytorch/pytorch
cd pytorch

# <!重要> 配置TCCL的安装路径
export USE_SYSTEM_NCCL=1
export NCCL_INCLUDE_DIR="/opt/tencent/tccl/include"
export NCCL_LIB_DIR="/opt/tencent/tccl/lib"

# 参考官网添加其他编译选项

# 安装开发环境
python setup.py develop
```

### 3、配置TCCL环境变量

```
export NCCL_DEBUG=INFO
export NCCL_SOCKET_IFNAME=eth0
export NCCL_IB_GID_INDEX=3
export NCCL_IB_DISABLE=0
export NCCL_IB_HCA=mlx5_bond_0,mlx5_bond_1,mlx5_bond_2,mlx5_bond_3,mlx5_bond_4,mlx5_bond_5,mlx5_bond_6,mlx5_bond_7
export NCCL_NET_GDR_LEVEL=2
export NCCL_IB_QPS_PER_CONNECTION=4
export NCCL_IB_TC=160
export NCCL_IB_TIMEOUT=22
export NCCL_PXN_DISABLE=0
export TCCL_TOPO_AFFINITY=4
```

注意：

需要通过TCCL\_TOPO\_AFFINITY=4开启网络拓扑感知特性。

### 4、验证 Pytorch

运行单机多卡或者多机多卡训练过程中有如下打印（export NCCL\_DEBUG=INFO）：

```
vm-3-17-centos:74350:74350 [0] NCCL INFO Bootstrap : Using eth0:10.100.3.17<0>
vm-3-17-centos:74350:74350 [0] NCCL INFO NET/Plugin : No plugin found (libnccl-net.so), using internal implementation
vm-3-17-centos:74350:74350 [0] NCCL INFO NET/IB : Using [0]mlx5_bond_0:1/RoCE [R0]; OOB eth0:10.100.3.17<0>
vm-3-17-centos:74350:74350 [0] NCCL INFO Using network IB
NCCL version 2.12.12_TCCL_v1.5+cuda11.6
vm-3-17-centos:74352:74352 [2] NCCL INFO Bootstrap : Using eth0:10.100.3.17<0>
vm-3-17-centos:74352:74352 [2] NCCL INFO NET/Plugin : No plugin found (libnccl-net.so), using internal implementation
vm-3-17-centos:74352:74352 [2] NCCL INFO NET/IB : Using [0]mlx5_bond_0:1/RoCE [R0]; OOB eth0:10.100.3.17<0>
vm-3-17-centos:74352:74352 [2] NCCL INFO Using network IB
```

## 5、验证nccl-tests

运行 nccl-tests 之前需要 export 对应的 TCCL 路径：

```
export LD_LIBRARY_PATH=/opt/tencent/tccl/lib:$LD_LIBRARY_PATH
```

## 6、软件版本支持

目前 TCCL 对应 NCCL 版本 2.12，要求 glibc 版本 2.17 以上，CUDA 版本 10.0 以上。其他 CUDA 版本支持请联系您的售前经理获取支持。

方法二：安装 Pytorch 通信插件

Pytorch 支持通过插件的方式接入第三方通信后端，所以在不重新编译 Pytorch 的前提下，用户可以使用 TCCL 通信后端，API 与 NCCL 完全兼容。详情可参考 [Pytorch 现有通信后端介绍](#)。

### 1、安装 Pytorch 通信插件

```
# 卸载现有的tccl和NCCL插件
dpkg -r tccl && dpkg -r nccl-rdma-sharp-plugins

# 卸载torch_tccl
pip uninstall -y torch-tccl

# 安装torch_tccl 0.0.2版本
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/tccl/torch_tccl-0.0.2_pt1.12-py3-none-any.whl && pip install torch_tccl-0.0.2_pt1.12-py3-none-any.whl && rm -f torch_tccl-0.0.2_pt1.12-py3-none-any.whl
```

### 2、修改业务代码

```
import torch_tccl

#args.dist_backend = "nccl"
args.dist_backend = "tccl"
torch.distributed.init_process_group(
    backend=args.dist_backend,
    init_method=args.dist_url,
```

```
world_size=args.world_size, rank=args.rank
)
```

### 3、配置TCCL环境变量

```
export NCCL_DEBUG=INFO
export NCCL_SOCKET_IFNAME=eth0
export NCCL_IB_GID_INDEX=3
export NCCL_IB_DISABLE=0
export NCCL_IB_HCA=mlx5_bond_0,mlx5_bond_1,mlx5_bond_2,mlx5_bond_3,mlx5_bond_4,mlx5_bond_5,mlx5_bond_6,mlx5_bond_7
export NCCL_NET_GDR_LEVEL=2
export NCCL_IB_QPS_PER_CONNECTION=4
export NCCL_IB_TC=160
export NCCL_IB_TIMEOUT=22
export NCCL_PXN_DISABLE=0
export TCCL_TOPO_AFFINITY=4
```

注意：

需要通过 `TCCL_TOPO_AFFINITY=4` 开启网络拓扑感知特性。

### 4、验证 Pytorch

您在执行分布式训练业务时，出现如下提示可确认通信后端被正确加载。

```
vm-3-17-centos:35915:35915 [0] NCCL INFO NET/Plugin : No plugin found (libnccl-net.so), using internal implementation
vm-3-17-centos:35915:35915 [0] NCCL INFO NET/IB : Using [0]mlx5_bond_0:1/RoCE [R0]; OOB eth0:10.100.3.17<0>
vm-3-17-centos:35915:35915 [0] NCCL INFO Using network IB
NCCL version 2.12.12_TCCl_v1.5+cuda11.6
vm-3-17-centos:35919:35919 [4] NCCL INFO Bootstrap : Using eth0:10.100.3.17<0>
vm-3-17-centos:35919:35919 [4] NCCL INFO NET/Plugin : No plugin found (libnccl-net.so), using internal implementation
vm-3-17-centos:35921:35921 [6] NCCL INFO Bootstrap : Using eth0:10.100.3.17<0>
vm-3-17-centos:35920:35920 [5] NCCL INFO Bootstrap : Using eth0:10.100.3.17<0>
```

### 5、软件版本限制

当前安装包仅支持Pytorch 1.12，其他 Pytorch 和 CUDA 版本支持请联系您的售前经理获取支持。

说明：

如果运行`nccl-tests`或者其他需要动态链接通信库的场景，请使用方法一安装 TCCL。

(推荐) 方法三：安装 NCCL 插件

如果您已经安装了 NCCL，也可以使用 NCCL 插件的方式使用 TCCL 加速能力。

#### 1、安装 NCCL 插件

以 Ubuntu 20.04 为例，您可以使用以下命令安装插件。

```
# 卸载现有的tccl和nccl插件
dpkg -r tccl && dpkg -r nccl-rdma-sharp-plugins

# 下载安装nccl 1.2插件
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/nccl/nccl-rdma-sharp-plugins_1.2_amd64.deb && dpkg -i nccl-rdma-sharp-plugins_1.2_amd64.deb

# 请确保集群内使用nccl插件版本一致，以下为nccl 1.0版本下载安装命令，推荐使用稳定性更优的nccl 1.2版本
# wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/nccl/nccl-rdma-sharp-plugins_1.0_amd64.deb && dpkg -i nccl-rdma-sharp-plugins_1.0_amd64.deb && rm -f nccl-rdma-sharp-plugins_1.0_amd64.deb
```

如果您使用 CentOS 或 TencentOS，参考以下步骤安装：

```
# 卸载现有的nccl插件
rpm -e nccl-rdma-sharp-plugins-1.0-1.x86_64

# 下载安装nccl 1.2插件
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/nccl/nccl-rdma-sharp-plugins-1.2-1.x86_64.rpm && rpm -ivh --nodeps --force nccl-rdma-sharp-plugins-1.2-1.x86_64.rpm

# 请确保集群内使用nccl插件版本一致，以下为nccl 1.0版本下载安装命令，推荐使用稳定性更优的nccl 1.2版本
# wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/nccl/nccl-rdma-sharp-plugins-1.0-1.x86_64.rpm && rpm -ivh --nodeps --force nccl-rdma-sharp-plugins-1.0-1.x86_64.rpm && rm -f nccl-rdma-sharp-plugins-1.0-1.x86_64.rpm
```

## 2、获取拓扑排序的 IP 列表

NCCL 插件不需要依赖文件可提供 bonding 口动态聚合和全局 hash 路由两种优化。如果需要支持网络拓扑的亲感知，用户可以通过排序的 IP 列表来实现。

IP 排序可以按照如下方式完成：

### a. 准备 IP 列表文件

VPC IP 地址通过 `ifconfig eth0` 获取，每行1个节点 IP，格式如下：

```
root@VM-125-10-tencentos:/workspace# cat ip_eth0.txt
172.16.177.28
172.16.176.11
172.16.177.25
172.16.177.12
```

### b. 执行排序

```
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/tccl/get_rdma_order_by_ip.sh && bash
```

```
get_rdma_order_by_ip.sh ip_eth0.txt
```

注意：

- 所有节点都安装了 curl 工具（比如对于 Ubuntu，通过 apt install curl 安装）。
- 执行脚本的节点可以 ssh 免密访问其他所有节点。

### c. 查看排序后的 IP 列表文件

```
root@VM-125-10-tencentos:/workspace# cat hostfile.txt
172.16.176.11
172.16.177.12
172.16.177.25
172.16.177.28
```

### 3、配置 TCCL 环境变量

```
export NCCL_DEBUG=INFO
export NCCL_SOCKET_IFNAME=eth0
export NCCL_IB_GID_INDEX=3
export NCCL_IB_DISABLE=0
export NCCL_IB_HCA=mlx5_bond_0,mlx5_bond_1,mlx5_bond_2,mlx5_bond_3,mlx5_bond_4,mlx5_bond_5,mlx5_bond_6,mlx5_bond_7
export NCCL_NET_GDR_LEVEL=2
export NCCL_IB_QPS_PER_CONNECTION=4
export NCCL_IB_TC=160
export NCCL_IB_TIMEOUT=22
export NCCL_PXN_DISABLE=0
```

```
# 机器 IP 手动排序之后，就不需要添加如下变量了
# export TCCL_TOPO_AFFINITY=4
```

### 4、修改启动脚本

您需要在启动分布式训练时修改启动脚本。例如，如果使用 deepspeed launcher 启动训练进程，拿到排序后的 IP 列表之后，将对应的 IP 列表写入 hostfile，再启动训练进程。

```
root@vm-3-17-centos:/workspace/ptm/gpt# cat hostfile
172.16.176.11 slots=8
172.16.177.12 slots=8
172.16.177.25 slots=8
172.16.177.28 slots=8

deepspeed --hostfile ./hostfile --master_addr 172.16.176.11 train.py
```

如果使用 torchrun 启动训练进程，通过 --node\_rank 指定对应的节点顺序，

```
// on 172.16.176.11
torchrun --nnodes=4 --nproc_per_node=8 --node_rank=0 --master_addr=172.16.176.11 train.py ...
// on 172.16.176.12
torchrun --nnodes=4 --nproc_per_node=8 --node_rank=1 --master_addr=172.16.176.11 train.py ...
// on 172.16.176.25
torchrun --nnodes=4 --nproc_per_node=8 --node_rank=2 --master_addr=172.16.176.11 train.py ...
// on 172.16.176.28
torchrun --nnodes=4 --nproc_per_node=8 --node_rank=3 --master_addr=172.16.176.11 train.py ...
```

如果使用 mpirun 启动训练进程，按照顺序排列 IP 即可。

```
mpirun \
-np 64 \
-H 172.16.176.11:8,172.16.177.12:8,172.16.177.25:8,172.16.177.28:8 \
--allow-run-as-root \
-bind-to none -map-by slot \
-x NCCL_DEBUG=INFO \
-x NCCL_IB_GID_INDEX=3 \
-x NCCL_IB_DISABLE=0 \
-x NCCL_SOCKET_IFNAME=eth0 \
-x NCCL_IB_HCA=mlx5_bond_0,mlx5_bond_1,mlx5_bond_2,mlx5_bond_3,mlx5_bond_4,mlx5_bond_5,ml
x5_bond_6,mlx5_bond_7 \
-x NCCL_NET_GDR_LEVEL=2 \
-x NCCL_IB_QPS_PER_CONNECTION=4 \
-x NCCL_IB_TC=160 \
-x NCCL_IB_TIMEOUT=22 \
-x NCCL_PXN_DISABLE=0 \
-x LD_LIBRARY_PATH -x PATH \
-mca coll_hcoll_enable 0 \
-mca pml ob1 \
-mca btl_tcp_if_include eth0 \
-mca btl ^openib \
all_reduce_perf -b 1G -e 1G -n 1000 -g 1
```

# 配置 HARP 分布式训练环境

## 操作场景

本文介绍如何通过云服务器控制台，为实例配置 HARP 分布式训练环境。

## 操作步骤

### 绑定弹性网卡

弹性网卡数量等于 GPU 卡的数量，例如8卡训练机器则需要绑定8张弹性网卡（加主网卡共9张网卡）。具体步骤如下：

1. 登录云服务器控制台，选择实例 ID 进入详情页面。
2. 在实例详情页中，选择弹性网卡页签，并单击绑定网卡。
3. 在弹出的绑定弹性网卡窗口中，选择弹性网卡，单击确认即可。

### 配置并验证环境

1. 参见使用标准登录方式登录 Linux 实例，登录实例。
2. 执行以下命令，执行配置脚本。

```
curl -s -L http://mirrors.tencent.com/install/GPU/taco/taco_setup.sh | sudo bash
```

返回结果如下：

```
[Thu 12 May 2022 08:36:50 PM CST] ===== Start configuring environment for TACO-Training =====
[Thu 12 May 2022 08:36:50 PM CST] Check whether assistant nics are attached and enough
[Thu 12 May 2022 08:36:50 PM CST] Modify /etc/rc.local to auto re-configure harp after each reboot
[Thu 12 May 2022 08:36:50 PM CST] Bind assistant nic(s) to kernel at first
[Thu 12 May 2022 08:36:51 PM CST] Remove existing harp configurations
[Thu 12 May 2022 08:36:51 PM CST] Start downloading tools package ...
[Thu 12 May 2022 08:36:51 PM CST] Start loading the uio module for "Ubuntu" ...
[Thu 12 May 2022 08:36:51 PM CST] Start generating harp config files ...
[Thu 12 May 2022 08:36:53 PM CST] Set up HARP successfully
```

3. 执行以下命令，重启实例。

```
sudo reboot
```

4. 依次执行以下命令，检查是否配置成功。

- 检查大页内存是否配置成功：

```
cat /proc/meminfo | grep HugePages_Total
```

返回如下结果，表示配置成功。

```
HugePages_Total: 50
```

- 检查是否产生了配置文件：

```
ls -l /usr/local/tfabric/tools/config/ztcp*.conf
```

返回结果如下，表示已产生配置文件。

```
(base) ubuntu@VM-18-223-ubuntu:~$ ls -l /usr/local/tfabric/tools/config/ztcp*.conf
-rw-rw-rw- 1 root root 898 May 25 10:36 /usr/local/tfabric/tools/config/ztcp1.conf
-rw-rw-rw- 1 root root 900 May 25 10:36 /usr/local/tfabric/tools/config/ztcp2.conf
-rw-rw-rw- 1 root root 901 May 25 10:36 /usr/local/tfabric/tools/config/ztcp3.conf
-rw-rw-rw- 1 root root 905 May 25 10:36 /usr/local/tfabric/tools/config/ztcp4.conf
-rw-rw-rw- 1 root root 906 May 25 10:36 /usr/local/tfabric/tools/config/ztcp5.conf
-rw-rw-rw- 1 root root 911 May 25 10:36 /usr/local/tfabric/tools/config/ztcp6.conf
-rw-rw-rw- 1 root root 912 May 25 10:36 /usr/local/tfabric/tools/config/ztcp7.conf
-rw-rw-rw- 1 root root 894 May 25 10:36 /usr/local/tfabric/tools/config/ztcp.conf
```

# 配置容器 SSH 免密访问

## 操作场景

本文介绍如何配置服务器间的容器 SSH 免密访问。

## 操作步骤

### 说明

本文以两台机器间容器 SSH 免密访问为例，配置步骤需在两台机器上同步。

1. 参见 [使用标准登录方式登录 Linux 实例](#)，登录实例。
2. 执行以下命令，允许 root 使用 ssh 服务，并启动服务（默认端口：22）。

```
sed -i 's/#PermitRootLogin prohibit-password/PermitRootLogin yes/' /etc/ssh/sshd_config
```

3. 依次执行以下命令，修改容器内 ssh 默认端口为2222，防止与 host 所使用的22端口冲突。

```
sed -i 's/#Port 22/Port 2222/' /etc/ssh/sshd_config
```

```
service ssh restart && netstat -tulpn
```

4. 执行以下命令，设置 root 密码。

```
passwd root
```

5. 执行以下命令，产生 SSH Key。

```
ssh-keygen
```

6. 创建 `~/.ssh/config`，并添加以下内容后，保存并退出，完成 host alias 配置。

```
# ！注意：  
# 如果是CVM机型，则ip是两台机器`ifconfig eth0`显示的ip
```

```
# 如果是RDMA机型, 则ip是两台机器`ifconfig bond0`显示的ip
Host gpu1
hostname 10.0.2.8
port 2222
Host gpu2
hostname 10.0.2.9
port 2222
```

7. 使用 `ssh-copy-id` 将 SSH key 拷贝到对应的机器, 使两台机器互相免密, 同时本机对自身进行免密。

```
ssh-copy-id gpu1
ssh-copy-id gpu2
```

# 容器安装用户态 RDMA 驱动

## 操作场景

本文介绍如何为容器安装用户态 RDMA 驱动。

## 操作步骤

### 说明

本文以 Ubuntu 20.04 操作系统的机器为例。

1. 执行以下命令，下载对应容器中的 OS 版本的 MLNX OFED 驱动。

```
wget https://www.mellanox.com/downloads/ofed/MLNX_OFED-5.4-3.1.0.0/MLNX_OFED_LINU  
X-5.4-3.1.0.0-ubuntu20.04-x86_64.tgz
```

若您使用了其他版本操作系统，则请访问 [Linux InfiniBand Drivers](#) 下载对应的版本。选择步骤如下：

### 注意

OFED 版本选择 5.4-3.1.0.0。

GPUDirect RDMA  
CloudX Software  
FlexBoot  
UEFI  
mlxup - Firmware Utility  
MFT - Firmware Tools  
IB Management &  
Monitoring Tools  
Common Information Model

**Linux Inbox Drivers**  
Linux Drivers for Ethernet and InfiniBand adapters are also available Inbox in all the major distributions, RHEL, SLES, Ubuntu and more.

View the explanation of [MLNX\\_OFED OS support models](#) and information about OFED LTS for NVIDIA products.

Benefits   Download   **LTS Download**

**Note:** MLNX\_OFED 4.9-x LTS should be used by customers who would like to utilize one of the following:

- NVIDIA ConnectX-3 Pro
- NVIDIA ConnectX-3
- NVIDIA Connect-IB
- RDMA experimental verbs library (mlnx\_lib)
- OSs based on kernel version lower than 3.10

**Note:** All of the above are not available on MLNX\_OFED 5.x branch.

**Note:** MLNX\_OFED 5.4-x LTS should be used by customers who would like to utilize NVIDIA ConnectX-4 onwards adapter cards and keep using stable 5.4-x deployment and get:

- Critical bug fixes
- Support for new major OSs

[Inbox Drivers - Documents](#)

**CONFIGURATION TOOLS**

**ACADEMY ONLINE COURSES**

**REQUEST A DEMO**

**MLNX\_OFED LTS Download Center**

Current Versions   Archive Versions   START OVER

Version (Archive)	OS Distribution	OS Distribution Version	Architecture	Download/Documentation
5.4-3.4.0.0	Ubuntu	Ubuntu 21.04	x86_64	ISO: <a href="#">MLNX_OFED_LINUX-5.4-3.1.0.0-ubuntu20.04-x86_64.iso</a>
5.4-3.2.7.2.3	SLES	Ubuntu 20.04	ppc64le	SHA5256: 56a99c45bc65f57ada11965d50900127968983978885e2c9d7c45a6410e8f6be
5.4-3.1.0.0	RHEL/CentOS	Ubuntu 18.04	aarch64	Size: 399M
5.4-3.0.3.0	Oracle Linux	Ubuntu 16.04		<b>tgz: <a href="#">MLNX_OFED_LINUX-5.4-3.1.0.0-ubuntu20.04-x86_64.tgz</a></b>
4.9-4.1.7.0	OPENEULER	Ubuntu 14.04		SHA5256: a8540ff82c050d103d5a499cb7671657278cf2dbfcf9cd553f43c51cb798eede
4.9-4.0.8.0	KYLIN			Size: 398M
4.9-4.0.8.0	Fedora			SOURCES: <a href="#">MLNX_OFED_SRC-debian-5.4-3.1.0.0.tgz</a>
4.9-4.0.8.0	EulerOS			
4.9-4.0.8.0	Debian			
4.9-4.0.8.0	Citrix			
4.9-4.0.8.0	XenServer			
4.9-4.0.8.0	Host			

2. 依次执行以下命令，进行解压及安装。

```
tar xf MLNX_OFED_LINUX-5.4-3.1.0.0-ubuntu20.04-x86_64.tgz
```

```
cd MLNX_OFED_LINUX-5.4-3.1.0.0-ubuntu20.04-x86_64
```

```
./mlnxofedinstall --user-space-only --without-fw-update --force
```

安装过程中出现的红色 warning 信息可忽略，直至页面出现 Installation passed successfully 绿色字样，表示安装成功。

## 相关操作

若您在安装过程中出现如下图所示错误：

```
Installing mlnx-iproute2-5.6.0...
Installing neohost-backend-1.5.0...
Failed to install neohost-backend DEB
Collecting debug info...
See /tmp/MLNX_OFED_LINUX.7204.logs/neohost-backend.debinstall.log
```

请参考以下步骤处理：

1. 由于 neohost 需要依赖 python2，执行以下命令，修改系统默认的 python 版本。

```
ln -sf /usr/bin/python2.7 /usr/bin/python
```

2. 执行以下命令，确认 python 版本。

```
python --version
```

如果提示找不到 python 命令，则需要安装 python2.7。

3. 执行以下命令，重新安装 ofed。

```
./mlnxofedinstall --user-space-only --without-fw-update --force
```

4. 执行以下命令，恢复 python3 作为默认 python 版本。

```
update-alternatives --install /usr/bin/python python /usr/bin/python3 1
```

# TACO LLM 推理加速引擎

## TACO-LLM 部署

### TACO-license 部署

TACO-license 是 TACO-LLM 项目中用于鉴权和记录日志的组件。

#### TACO-license 背景及适用范围

TACO-license 是 TACO-LLM 项目下用于鉴权和记录日志的组件，适用于私有化部署 TACO-LLM。私有化客户需要在集群中部署 server 端，以对 client 端（即 TACO-LLM 本身）的请求进行鉴权。搭建条件如下，本文档中将使用 `$your_license_version` 来代表实际的版本号，请在代码中进行相应替换：

```
//搭建环境：需要K8s，server以service的形式部署并提供服务  
docker images: taco-license-server:$your_license_version
```

成功鉴权后，可以正常使用 TACO-LLM 功能。

#### 签发 license

只有被授权并在过期时间之前，在最大同时访问数量（如果设定）的限制下，才能成功鉴权的 GPU 才能使用。最大同时访问数量指的是同时访问 TACO-license 的 GPU 的最大数量，例如：授权10000台 GPU，同时限制最大同时访问的 GPU 数量为10台。

1. 客户将需要签发 license 的 GPU 发送给 TACO 团队，其方法为在所有可能执行 GPU 的机器上执行：

```
nvidia-smi -L
```

将结果写入文件（例如 `request.txt`），多台机器的结果直接按顺序粘贴合并到一个文件中，并提供过期时间以及最大访问数量，交给 TACO 团队以签发许可证。

2. 签发的许可证文件名为 `license-taco_xxxxxxx.dat`（例如 `license-taco_1234-leolingli-20240906104354.dat`）。

#### 搭建 taco-license-server

1. 准备好 `taco-license-server` 的镜像，可以下载到本地或者上传到可供下载的位置。此包由 TACO 团队提供，下载地址为：

```
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/taco-llm/license-server/latest/taco-license-server-latest.tgz
```

```
#将docker image 解压或放到集群可以下载到的地址，如：  
docker load -i taco-license-server-latest.tgz
```

```
#docker pull or docker tag taco-license-server: $your_license_version \#$your_license_repository/taco-license-server:$your_license_version
```

2. 使用 Helm Chart 搭建服务，可以参考以下示例：

```
#install helm first
#wget taco-operator-idc package
wget https://taco-1251783334.cos.ap-shanghai.myqcloud.com/taco-llm/license-server/latest/taco-operator-idc-latest.tgz

#假设您的taco-license-server地址为 $your_license_repository/taco-license-server:$your_license_version
helm install --generate-name \
  --set global.repository="$your_license_repository" \

./taco-operator-idc-latest.tgz
```

3. 服务将部署在 kube-system 命名空间中，可以通过 curl 命令来验证是否成功搭建，截图如下：

```
./taco-operator-idc-v0.1.0.tgz
[root@vm-1-143-centos ~]# k get svc
NAME                TYPE          CLUSTER-IP    EXTERNAL-IP    PORT(S)          AGE
kubernetes          ClusterIP    10.10.0.1     <none>         443/TCP          87d
taco-license-server ClusterIP    10.10.195.83 <none>         10080/TCP        20h
[root@vm-1-143-centos ~]# curl 10.10.195.83:10080
taco license server
[root@vm-1-143-centos ~]#
```

## 导入 license

将上述签发的 license ( license-taco\_xxxxxxx.dat ) 导入集群，使 license 生效：

1. 在集群中执行以下操作以导入许可证。请注意，导入操作将覆盖所有现有的许可证。因此，新签发的许可证需要替代之前的所有许可证：

```
LICENSE_FILE_NAME=license-taco_xxxxxxxxxx.dat
```

```
LICENSE_SERVER_IP=$(kubectl get svc -n kube-system taco-license-server -o jsonpath="{.spec.clusterIP}")
```

```
curl -X POST -F "file=@$LICENSE_FILE_NAME" http://${LICENSE_SERVER_IP}:10080/ls
```

2. 执行上述命令后，可以看到提示，显示 license 已成功导入。

```
[root@vm-1-143-centos install]# curl -X POST -F "file=@$LICENSE_FILE_NAME" http://${LICENSE_SERVER_IP}:10080/ls
license file sucessfully uploaded
[root@vm-1-143-centos install]#
```

您也可以通过以下方法查看签发的许可证信息：

```
curl your-service-ip:10080/ls/view
```

```
license file sucessfully uploaded
[root@vm-1-143-centos install]# curl 10.10.195.83:10080/ls/view
appid: "1234"
customer: leolingli
accesslimit: 0
gpuinfos:
- uuid: GPU-4c34f7a2-cd33-469b-019f-9af0bb76a073
  brand: Tesla T4
  issuedtime: 2024-09-06T10:43:54.427014401+08:00
  expiredtime: 2025-09-06T10:43:54.427014401+08:00
- uuid: GPU-47326f70-8b22-8137-0b48-1d3746ff7581
  brand: Tesla T4
  issuedtime: 2024-09-06T10:43:54.427014401+08:00
  expiredtime: 2025-09-06T10:43:54.427014401+08:00
- uuid: GPU-3adbf791-a221-e4cd-2ab0-e14f3ed47e19
  brand: Tesla T4
  issuedtime: 2024-09-06T10:43:54.427014401+08:00
  expiredtime: 2025-09-06T10:43:54.427014401+08:00
- uuid: GPU-88bbd1c8-4ac5-ee9e-a1b1-5caf134be771
  brand: Tesla T4
  issuedtime: 2024-09-06T10:43:54.427014401+08:00
  expiredtime: 2025-09-06T10:43:54.427014401+08:00
- uuid: "1234"
  brand: Fake Brand
  issuedtime: 2024-09-06T10:43:54.427014401+08:00
  expiredtime: 2025-09-06T10:43:54.427014401+08:00
```

client 端使用

1. 在容器中，只需使用 `taco-llm` 即可，无需特殊配置。您可以尝试在 `pod` 中执行以下操作：

```
curl http://taco-license-server.kube-system.svc.cluster.local:10080
```

2. 可以看到相同的输出：

```
taco license server
```

这意味着客户端可以访问许可证。

特殊情况：自定义 `license` 地址

说明：

本章节适用于使用非 `kube-system` 命名空间或希望自定义 `qgpu-license-server` 部署方式的用户。

`Taco-license-server` 默认支持两种连接方式，即通过 <http://taco-license-server.kube-system.svc.cluster.local:10080> 和 <http://taco-license-server:10080> 访问。这意味着部署在 `kube-system` 命名空间和与负载 `pod` 处于同一命名空间的服务地址都可以进行访问。如果用户想在其他命名空间使用，或者由于配置原因无法公开 `taco-license-server` 地址以允许客户端鉴权，则需要用户自行解决 `taco-license-server` 的部署问题。这包括部署一个可以访问到 `taco-license-server` 地址的 `pod`，并且设置 `TACO_LS_ADDR` 环境变量为 `taco-license-server` 的地址。

例如，如果 `taco-license-server` 可以通过地址 <http://10.10.161.72:10080> 访问，则应在 `pod` 的环境变量中添加以下变量：

- 可以通过配置 `YAML` 文件中的环境变量来实现。

```
[root@vm-1-143-centos ~]# cat taco-0.yaml
apiVersion: v1
kind: Pod
metadata:
  name: taco-0
spec:
  containers:
  - name: taco-0
    image: peaceforever/taco/taco-infer:v99.0.2
    env:
    - name: TACO_LS_ADDR
      value: "http://10.10.161.72:10080"
    command: ["sleep", "12345000"]
```

- 或在容器中 bashrc 文件中添加变量：

```
export TACO_LS_ADDR=http://$your_license_addr:10080
```

Taco-LLM 将尝试通过此地址访问 Taco-License-Server。

#### 注意：

如果要显式设置 TACO\_LS\_ADDR 的值，请确保将 TACO\_LS\_ADDR 的值导出为正确的 taco-license-server 地址。尝试访问无效地址会导致鉴权程序等待超时（通常为30秒），从而延迟鉴权启动时间，可能会影响到最初几次鉴权的结果。

# TACO LLM 安装

## 环境准备

TACO-LLM 需要依赖 GPU 相关的基础软件，如 GPU 驱动 /CUDA 等。为了避免基础软件依赖导致 TACO-LLM 无法正常运行，我们提供了 TACO-LLM docker环境镜像，建议您优先使用该镜像作为 TACO-LLM 的运行环境。按照如下命令可以获取 docker 镜像并启动容器环境：

```
docker run -it \  
  --privileged \  
  --net=host \  
  --ipc=host \  
  --shm-size=16g \  
  --name=taco_llm \  
  --gpus all \  
  -v /home/workspace:/home/workspace \  
  ccr.ccs.tencentyun.com/taco/tacollm-dev:latest /bin/bash
```

## 安装 whl 包

1. 按照如下命令在容器环境中安装 TACO-LLM：

```
pip3 install taco_llm-${version}-cp310-cp310-linux_x86_64.whl
```

2. 安装 TACO-LLM whl 包时，会自动安装相关的 python 依赖包。

# TACO LLM 使用

## 离线模式

在离线推理业务场景中，您可以通过离线模式使用 TACO-LLM。本文档通过一个简单的例子介绍了如何使用 TACO-LLM 的离线模式。

## 导入 LLM 和 SamplingParams

首先，需要从 `taco_llm` 导入所需使用的 LLM 和 `SamplingParams` 类：

```
from taco_llm import LLM, SamplingParams
```

## 构建 prompts 和采样参数

接下来，构建所需的 prompts 和采样参数。本示例构建了4条 prompt，并设置了采样参数，其中 `temperature` 为 0.8，`top_p` 为0.95。完整的采样参数配置可以参见[采样参数 API](#)。

```
# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)
```

## 构建 LLM 对象

接下来，我们将构建 LLM 实例。本示例使用了 `facebook/opt-125m` 模型以及其他默认配置参数来构建 LLM 实例。完整的 LLM 构建参数配置可以参见 [离线API](#)。

```
# Create an LLM.
llm = LLM(model="facebook/opt-125m")
```

## 推理计算

最后，调用 LLM 对象的 generate 接口进行推理计算：

```
# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```

至此，TACO-LLM 的离线模式使用已经完成。以下是本示例的完整代码：

```
from taco_llm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
# Create a sampling params object.
sampling_params = SamplingParams(temperature=0.8, top_p=0.95)

# Create an LLM.
llm = LLM(model="facebook/opt-125m")
# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)
# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")
```

# 在线模式

TACO-LLM 提供了实现 OpenAI [Completions](#) 和 [Chat API](#) 的 HTTP 服务端，您可以按照以下流程进行使用。

## 启动服务

首先，执行以下命令启动服务：

```
taco_llm serve facebook/opt-125m --api-key taco-llm-test
```

## 发送请求

您可以使用 OpenAI 的官方 Python 客户端来发送请求：

```
from openai import OpenAI

client = OpenAI(
    base_url="http://localhost:8000/v1",
    api_key="taco-llm-test",
)

completion = client.chat.completions.create(
    model="facebook/opt-125m",
    messages=[
        {"role": "user", "content": "Hello!"}
    ]
)

print(completion.choices[0].message)
```

您也可以直接使用 HTTP 客户端来发送请求：

```
import requests

api_key = "taco-llm-test"

headers = {
    "Authorization": f"Bearer {api_key}"
}

payload = {
```

```
"prompt": "Hello!",
"stream": True,
"max_tokens": 128,
}

response = requests.post("http://localhost:8000/v1/completions",
                        headers=headers,
                        json=pload,
                        stream=True)

for chunk in response.iter_lines(chunk_size=8192,
                                decode_unicode=False,
                                delimiter=b"\0"):
    if chunk:
        data = json.loads(chunk.decode("utf-8"))
        output = data["text"][0]
        print(output)
```

## 完整服务端参数配置

执行 `taco_llm serve -h` 命令可以查看 TACO-LLM 完整的在线模式参数配置，详细内容请参见：[在线模式 API](#)。

## 完整客户端参数配置

除了少部分参数不支持外，TACO-LLM 完全支持 OpenAI 的参数配置。您可以参见 [OpenAI API 官方文档](#) 查看完整的 API 参数配置。不支持的少部分参数配置如下：

- Chat: tools, and tool\_choice。
- Completions: suffix。

# 基础配置

## 基本配置

以一个实际运行的例子进行离线测试为例：此例子可以参考离线模型的编写方式。

```
python3 offline_test.py \  
  --model /models/Llama-2-7b-chat-hf/ \  
  --tokenizer /models/Llama-2-7b-chat-hf/ \  
  --dataset /datasets/ShareGPT_V3_unfiltered_cleaned_split.json \  
  --num-prompts 16 \  
  --max-num-batched-tokens 10240 \  
  --max-num-seqs 32 \  
  --trust-remote-code
```

参数说明：

- --model：指定模型的存放路径，支持相对路径、绝对路径或者 repo\_id 则会从 huggingface 上下载，建议使用本地路径可加快运行速度。
- --tokenizer：同模型同一个参数。
- --dataset：指定当前测试 case 需要从哪个数据集采集 prompt 数据。
- --num-prompts：表示请求数。
- --max-num-batched-tokens：表示每次执行推理时支持最长的处理 token 数，多个 batch 的总和数，如果每个请求长，则会分多批完成请求。
- --max-num-seqs：后端支持最大的 batch 数，配置值建议参考实际 GPU 总显存及处理请求的长短值考虑。如果配置过大启动服务会慢，配置过小可能影响最终的吞吐性能。
- --trust-remote-code：表示信任当前部署，如果使用 model repo\_id 时不加此参数可能导致模型加载失败。除了上面参数外，还有一些可选配置参数供参考。

## 多卡推理

```
--tensor-parallel-size 1 \  

```

--tensor-parallel-size：tensor 并行，支持1, 2, 4, 8，会被模型的 layer 层数整除即可。

## cudaGraph 优化

```
--enforce-eager \
```

**\*\*--enforce-eager :** \*\*默认支持 CUDA Graph 优化。如果不想使用 CUDA Graph , 可以添加 --enforce-eager 参数。默认配置的 CUDA Graph 在启动模型时会捕获图的时间开销, max-num-seqs 配置越大, 所需时间越长。

## 显存占用配置

```
--gpu-memory-utilization 0.9 \
```

```
--conservative-dry-run \
```

- --gpu-memory-utilization : 小数值取值范围 : 0.1 ~ 0.95 表示 GPU 卡占用显存, 这部分显存主要有三部分: 进程初始化显存 (非 torch)、权重、kvcache; 剩下的显存一般用于非 torch 分配、cuda graph 等场景, 如果此值配置过高, max-num-seqs 配置过大, 会导致 cuda graph 显存不够, 导致初始化失败。
- --conservative-dry-run : 增加此参数表示在量化等场景由于内部 kernel 也使用显存, 而这部分显存并未统计到上述利用值上面, 导致 OOM, 建议在量化场景下增加此参数保证安全运行。

## speculative推理

```
--speculative-model /models/*** --num-speculative-tokens 3 \
```

- --speculative-model : 后面跟 speculative 模型的路径, 运行绝对路径、相对路径等。
- --num-speculative-tokens : 3 表示 speculative 模型一次生成的 token 数, 可以自行配置调优。

## 多步推理

```
--num-scheduler-steps 4 \
```

--num-scheduler-steps : 支持多步处理, 可以配置, [1-8]等整数值, 在 decoding 阶段, 部分 cpu 结果的处理可以 overlap 到 GPU 上, 增加吞吐值, 默认不打开。

# Lookahead Cache

## 默认方式开启

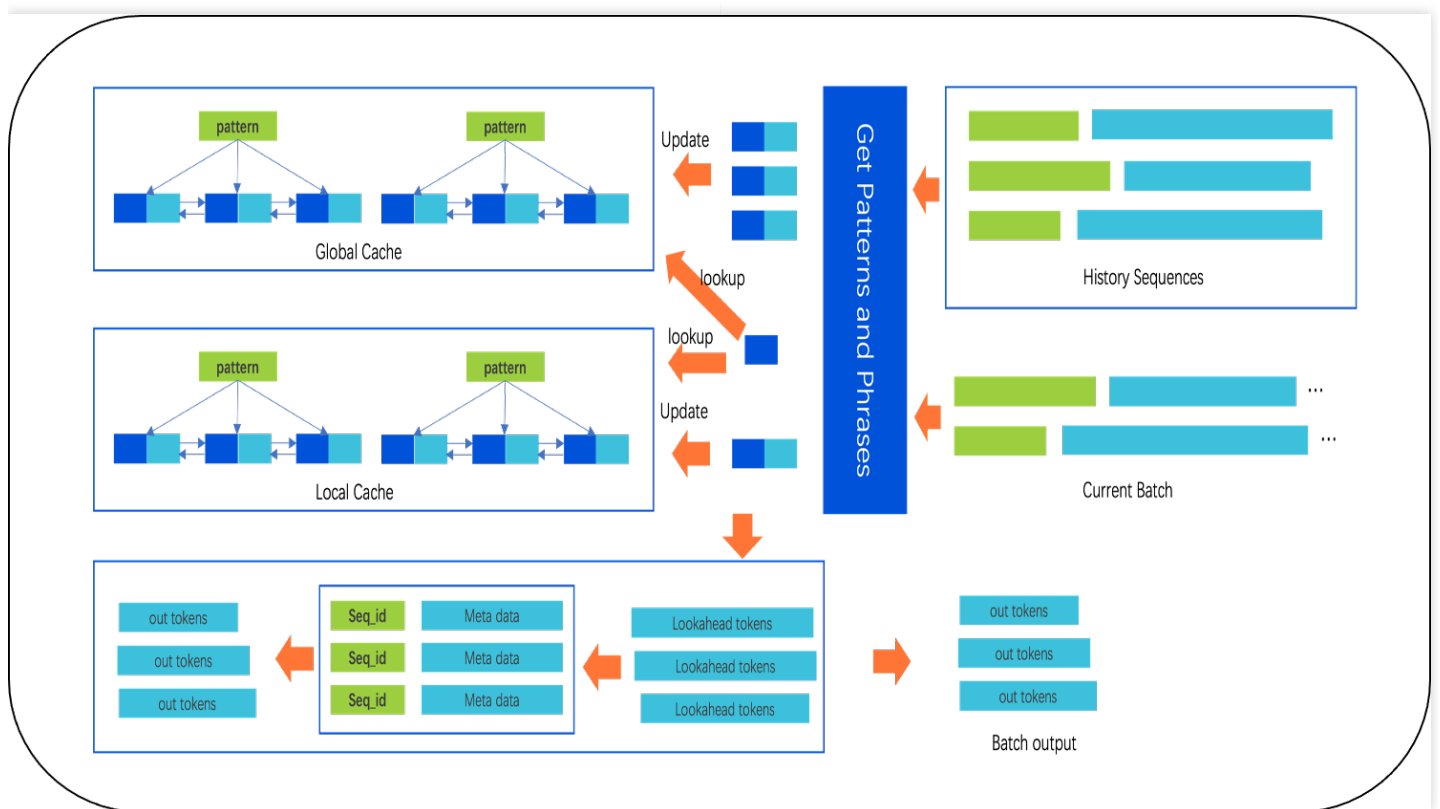
您只需在启动命令行中加入以下命令，即可开启 lookahead-cache：

```
--lookahead-cache-config-dir ./ # 或者任意其他目录
```

如果您只想了解 LookaheadCache 在默认配置下的加速效果（默认配置通常已足够优秀），那么以上便是您需要知道的所有信息。如果您希望进一步调整性能，选择不同的 CacheMode，或查看 debug log 等，那么您需要继续阅读后续内容。

## 基本原理

一句话概括 Lookahead Cache 的基本原理：使用历史的 Token 对（key: value）来预测当前 Token 对的 value 值。所提及的 Lookahead Cache 是 Taco-LLM 的 Lookahead 技术，主要设计并实现了两个基础方案。下图展示了第一个方案的 Lookahead Cache：



- 方案一的 Lookahead Cache 参见了LLMA 和 N-gram 的设计思想，并在此基础上进一步优化，并增加了新特性。

主要特点是：

- 支持变长的 Lookahead Len；
  - 支持 batch 化；
  - 支持变长 cuda-graph。
- 方案二抛弃了 N-gram 和 LLMA 的设计思想，重新设计，相比方案一继承了其所有优点，同时拥有更加高的平均命中长度，但是也更加复杂，这里不详细叙述。

在功能上，方案二新增支持：

- Multi-path 的多候选结果输出；
- 基于前缀树的多候选结果的输出；
- 树的动态裁剪；整体上的预测结果命中率是方案一的1.4x，部分场景是方案一的2.0x。  
但是相比方案一，方案二的冷启动问题更加明显，所以往往结合方案一一起使用。

## 进阶使用方式

### 配置介绍

在这种使用方式下，您需要提供一个命名为 `lookahead_cache_config.json` 的配置文件，该文件必须位于通过 `--lookahead-cache-config-dir` 指定的目录中。（文件名：`lookahead_cache_config.json`）

```

{
  "cache_mode": 2,
  "cache_size": 5000000,
  "copy_length": 7,
  "match_length": 2,
  "turbo_match_length": 7,
  "min_match_length": 2,
  "cell_max_size": 16,
  "voc_size": 200000,
  "max_seq_len": 32768,
  "eos_token_id": 2,
  "top_k": 1,
  "threshold": 2.0,
  "decay": 0.5,
  "is_hybrid": true,
  "is_debug": true,
  "log_interval": 3000,
  "target_parallelism": 512,
  "top_k_in_cell": 16,
  "token_paths_top_k": 2,
  "start_freq": 10.0,
  "num_threads": 8,
  "global_cache_switch": true
}

```

还有选择性的 warmup\_file 配置。

参数	参数解释
cache_mode	0表示 RawLookaheadCache (第一版), 1表示 TurboLookaheadCache (第二版), 2表示两者混合的方式, 混合方式: 在 TurboLookaheadCache 返回结果为空时, 使用 RawLookaheadCache 来预测, RawLookahead 的命中率更低, 但是在开始阶段的触发概率更高, 而 TurboLookaheadCache 的 warmup 的时间稍长, 在测试数据很少的时候, 建议 cache_mode 设置为2, 如果有充足的测试数据, 建议设置 cache_mode 为1, 因为其命中率和准确率在大部分场景下都更好, 有更少的冗余计算。
cache_size	cache 的大小, 可以不用设置, 而通过 taco-llm 的命令行参数 --cpu-decoding-memory-utilization 来设置。它是一个0到1的值, 表示使用当前环境下内存的比值, 默认是0.15。
copy_length	lookahead 往前看的长度。
match_length	触发 lookahead cache 的匹配长度。
turbo_match_length	TurboLookaheadCache 的最大匹配长度。

参数	参数解释
min_match_length	TurboLookaheadCache 的最小匹配长度。
cell_max_size	二级缓存的大小，同一个匹配下不同词组，按照出现频率和时间，组成一个 LRU 的 cache。
voc_size	tokenizer 中词表的大小。
max_seq_len	模型的最大序列长度，不用配置，直接可以从模型的 config 中获取。
eos_token_id	结束 token，不用配置，直接从模型的 config 文件中获取。
top_k	取二级 cache 频率最大，时间最近的词的数目。
threshold	二级 cache 满了以后，需要删除时，要删除对象的频率阈值。
decay	对于二级 cache 满了时，且所有的频率都高于 threshold 时，需要先衰减。
is_hybrid	针对 RawLookaheadCache 打开，多种 match-length 会结合，例如，match-length=3 时，会混合 match-length=3, 2, 1。
is_debug	打开时，会输出 cache log。
log_interval	打印 log 的频次。
target_parallelism	最大并行度，即 $\sum(\text{seq\_lens})$ ，开启了就会，在 $\sum(\text{seq\_lens}) > \text{target\_parallelism}$ 后会根据各个 seq 的命中情况来惩罚 copy_length。仅仅对 cache_mode=0 时有效。
top_k_in_cell	仅针对 TurboLookaheadCache，表示查找二级缓存时，一次返回 token 的数目。
token_paths_top_k	表示 beam search 的宽度，默认是 1，在命中率比较低时，可以尝试把这个参数改为 2，或者更大，建议在 4 以内。这个参数的开启在一定程度上增加了命中率，同时也会增加冗余计算。
start_freq	表示 local cache 的初始比重，越高表示优先保留 Local cache 中得到的 path，仅在 token_paths_top_k > 1 时生效。
num_threads	表示在使用 TurboLookaheadCache 的并发度，提升的是 TurboLookaheadCache 本身的运行速度。
global_cache_switch	为 True 表示开启 global_cache，false 表示关闭，这样样本之间将不受影响。
ignore_prompt	忽略 prompt，多针对翻译等 prompt 和 generation tokens 毫无关系的场景。

## 注意的问题

ignore\_eos

正确操作：ignore\_eos=False

在测试 taco\_llm lookahead 时，这个参数需要设置为 False，因为 LookaheadCache 里，遇到 eos 的 token\_id 时，会把 local\_cache 清除掉。将会影响 hit\_rate 和 global\_average\_hit\_len 等指标，最终影响整体的性能。如果需要测试固定的输出长度，建议在搜索样本时，选择更长的输出样本，然后把 output\_len 固定，这样可以减少 output tokens 未到目标输出长度时，就遇到 eos token。

## 常用性能调节方式

注意：

以下的配置方式可以叠加使用。

### 小数据集测试

如果发现使用少量样本测试，TurboLookaheadCache (cache\_mode=1) 效果比 RawLookaheadCache (cache\_mode=0) 的性能还要差，此时建议用 cache\_mode=2 的混合模式。增加如下配置：

```
{
  "cache_mode": 2
}
```

这是由于 TurboLookaheadCache 的冷启动问题导致的。

### 性能不符合预期

现象：使用 lookahead-cache 默认配置会有 1.7x-3.x+ 的性能收益，如果没有达到这个收益。首先需要查看是否使用了 greedy 的采样方式（这样会有最好的性能）。如果因为业务的原因不能使用 greedy，可以在满足业务需要的情况下调小 temperature，尽可能的减少随机性。这样调节和 lookahead cache 的原理密切相关。lookahead cache 会将历史输入和输出的 tokens 都存放到 cache 中，如果输出的随机性太强，那么历史输入对当前输出的参考性，就会变弱，从而导致 cache 的命中率下降，加速性能也下降。

在上述配置都没有问题的情况下，依然不符合加速预期，可以在 lookahead cache 的配置目录下，增加含有如下内容的配置文件（文件名：lookahead\_cache\_config.json）：

```
{
  "is_debug": true,
  "log_interval": 300
}
```

添加此配置后，系统将输出 hit-rate 的相关指标，如下图所示。

```
INFO 04-28 16:10:04 lookahead.py:265] num_iters: 10785.0000, hit_iters:8836.4450, hit_len: 16887.0650, invalid_len: 3936.2950, hit_rate: 0.8249, global_average_hit_len: 1.5883, valid_average_hit_len: 1.9117, hit_valid_rate: 0.8135
INFO 04-28 16:10:27 lookahead.py:265] num_iters: 10800.0000, hit_iters:8848.6950, hit_len: 16915.2350, invalid_len: 3945.4950, hit_rate: 0.8250, global_average_hit_len: 1.5892, valid_average_hit_len: 1.9126, hit_valid_rate: 0.8136
INFO 04-28 16:10:51 lookahead.py:265] num_iters: 10815.0000, hit_iters:8860.6800, hit_len: 16942.6550, invalid_len: 3955.4150, hit_rate: 0.8252, global_average_hit_len: 1.5900, valid_average_hit_len: 1.9133, hit_valid_rate: 0.8137
INFO 04-28 16:11:16 lookahead.py:265] num_iters: 10830.0000, hit_iters:8872.8150, hit_len: 16970.5000, invalid_len: 3964.8750, hit_rate: 0.8253, global_average_hit_len: 1.5908, valid_average_hit_len: 1.9141, hit_valid_rate: 0.8137
```

## 各命中率指标的含义

- num\_iters: 平均每一个请求的 generation 阶段的自回归次数。
- hit\_iters : 平均每一个请求的命中长度大于0的迭代次数。
- hit\_len : 平均每一个请求的命中总长度。
- invalid\_len: 平均每一个请求 lookahead 的总长度减去命中的总长度: lookahead\_len - hit\_len。
- hit\_rate : 平均每一个请求命中的迭代次数除以总自回归次数: hit\_iters / num\_iters。
- global\_average\_hit\_len: 平均每一个请求的平均每次迭代的命中长度: hit\_len / num\_iters。
- valid\_average\_hit\_len: 平均每一个请求命中长度大于0的情况下, 平均命中长度 : hit\_len / hit\_iters。
- hit\_valid\_rate: 平均每一个请求在总的命中长度除以总的 lookahead 长度 : hit\_len/ lookahead\_len。

其中 global\_average\_hit\_len 加1, 即表示开启 lookahead 下, 每一次 decoding 迭代平均吐出的 token 数目, 理想情况下, 这个数据表示 decoding 过程的加速倍数。例如上图所示, global\_average\_hit\_len = 1.59, 加1为2.59, 那么理想情况下 decoding 阶段加速2.59倍。hit\_rate 表示至少有一个命中的命中次数与全局迭代次数的比值。valid\_average\_hit\_len 表示如果至少有一个 token 命中的情况下。平均的命中长度, 这个数据加1, 可以用来调节 copy\_length 的大小, 注意: 如果使用了 mutli-path 优化, 不能使用这个方式来调节。下面举例说明在采样方式没有问题时, 如何调节 Lookahead Cache。

## 指标 global\_average\_hit\_len < 0.8

当 global\_average\_hit\_len < 0.8 时, 属于命中长度较低的情况。首先, 确认是否开启了 MultiPath 功能。如果使用的是默认配置, 即没有指定配置文件或者 cache\_mode 项未进行配置, 那么可以确定 MultiPath 已开启, 且其值为 2。如果未开启, 可添加以下配置:

```
{
  "is_debug": true,
  "log_interval": 3000,
  "token_paths_top_k": 2,
  "start_freq": 10.0
}
```

- 添加该配置后, 会禁用并行度惩罚。因此, 在大批量处理 (>=64) 时, 冗余计算可能会显著增加, 这可能会影响性能。如果性能下降是由此原因引起的, 可以尝试降低 copy\_length 的值。如果这些操作都无效, 请检查 ignore\_eos 是否设置为 False (如上文所述)。
- 确定 MultiPath 生效, 但是 global\_average\_hit\_len 依然较低, 可以尝试增加路径的数量, 即调整 token\_paths\_top\_k 的设置。token\_paths\_top\_k 的设定可以采用 copy\_length / valid\_average\_hit\_len 的方法, 其中默认的 copy\_length 为7。
- 如果此时 global\_average\_hit\_len 指标仍然不高 (<0.8), 则需要检查测试场景, 例如生成序列长度很短 (<=32), 或者输入和输出为语音、多模态等毫无关联的场景。在这种情况下, Lookahead Cache 在原理上受限, 不会有很好的效果。

## 指标 global\_average\_hit\_len >= 0.6 && <= 1 性能差

global\_average\_hit\_len 在 [0.6, 1] 这个区间内，但是端到端的速度甚至没有提升，可以尝试以下操作：

- 大 batchsize( $\geq 32$ )，此时冗余计算会比较多，lookahead 带来的收益可能不足以抵消冗余计算带来的开销，可以试着调低 copy\_length，可以4，5，6分别尝试，如果测试数据足够，也是尝试将 cache\_mode 修改为 1，将 copy\_length 调节到2，3等。
- 低算力卡(H20, L20): 分析和调节方法与大 batchsize 场景类似，此时该问题可能会来的更早，例如 bs=16左右。

### 不同的样本请求差别较大时

不同的样本之间差别较大时，需要快速的更新 cache，增加时间局部性，可以尝试增加下面的一组调节参数。也可以通过 cell\_max\_size 指标的调节来改变cache 的变化速度，其值越小变化的越快，默认值是16，建议的调节范围是[8, 32].

```
{
  "start_freq": 1.0,
  "decay": 0.1,
  "cell_max_size": 16
}
```

### warmup

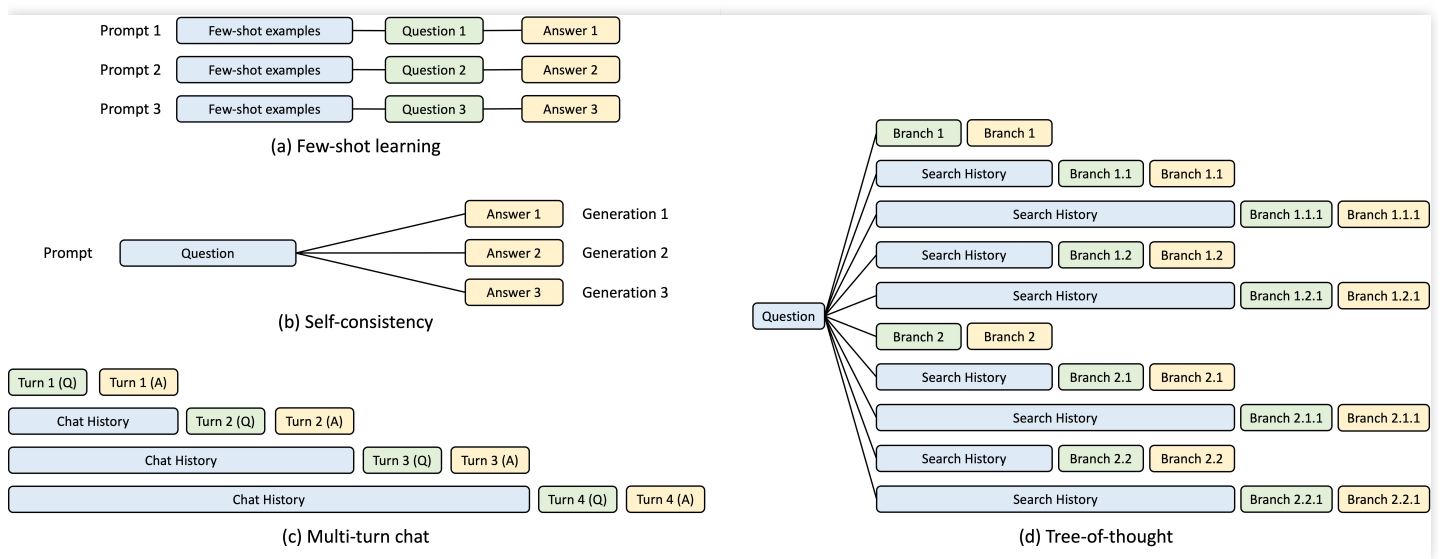
文件格式为 JSON 格式，数据格式为 [{"prompt": [token\_ids], "output": [token\_ids]}]。处理后，将其路径设置到 lookahead\_cache\_config.json 中，字段名为 "warmup\_file" : warmup\_file\_path。

# Auto Prefix Caching

在多线程、系统提示以及其他具有大量共同前缀的应用场景中，自动前缀缓存功能能够存储之前的前缀键值缓存，加快后续具有相同前缀请求的预填充阶段，从而降低首次响应时间，优化处理能力。

## 原理

LLM 推理计算主要分为两个过程：**Prefill 阶段 (Prompt 计算)**和 **Decode 阶段**。这两个阶段的计算特性存在不同，Prefill 阶段是计算受限的，而 Decode 阶段是访存受限的。为了避免重复计算，Prefill 阶段主要作用就是给 Decode 阶段准备 KV Cache。但这些 KV Cache 通常只是为单条推理请求服务的，当请求结束，对应的 KV-Cache 就会清除。那很自然的一种想法就是，KV Cache 能不能跨请求复用？在某些 LLM 业务场景下，多次请求的 Prompt 可能会共享同一个前缀 (Prefix)，比如少量样本学习，多线程等。在这些情况下，很多请求 Prompt 的前缀的 KV Cache 计算的结果是相同的，可以被缓存起来，给之后的请求复用。TACO-LLM 的 Auto Prefix Cache 技术可以针对这种场景进行优化，使得具有相同 Prompt 前缀的 KV-Cache 可以跨请求复用，降低计算开销，提升推理计算性能。



具体详情请参见：[Fast and Expressive LLM Inference with RadixAttention and SGLang](#)。

## 启动选项

```
--enable-prefix-caching
    Enables automatic prefix caching.
```

在 server 启动指令中添加该指令即可开启 Auto Prefix Caching 功能。

# Prefix Cache Offload

显卡的显存有限，除装载模型权重和运行激活空间以外，留给 prefix kv cache 的 blocks 数量是固定而有限的。当不同的 prefix 请求较多，随着请求的不断输入，之前的 prefix cache 就会被驱逐，之后有相同 prefix 的请求将无法命中 prefix cache。这导致需要重新执行 prefill 流程，从而无法获得加速效果。

TACO LLM 提供 prefix cache offload 功能，在显存 prefix cache 被驱逐时 offload 到 cpu 内存上，命中时 load 回 gpu 中，进而加速之前 prefix cache 被驱逐而得不到加速的请求。

## Offload 选项

- `--enable-prefix-cache-offload`  
Enables prefix cache offloading
- `--cpu-prefill-memory-utilization CPU_PREFILL_MEMORY_UTILIZATION`  
the memory is used for prefill cache, which can range from 0 to 1. If unspecified, will use the default value of 0.3.
- `--apc-offload-not-lazy`  
If set, lazy launch of layer 2~n-1 will be disabled.
- `--apc-offload-min-access-threshold APC_OFFLOAD_MIN_ACCESS_THRESHOLD`  
Min threshold for evict offloading. Default 1.
- `--apc-offload-enable-hit-cnt`  
Enable hit count in APC.
- `--apc-offload-gpu-evictor-limit APC_OFFLOAD_GPU_EVICTOR_LIMIT`  
The free table size limited in gpu evictor. -1 default

- `--enable-prefix-cache-offload`
  - 在开启了 APC 的基础上打开该开关，即可启用 prefix cache offload 功能。
- `--cpu-prefill-memory-utilization`
  - 用于 kv cache 的 cpu 内存比例，默认0.3，按卡均分。
  - Note：该比例按机器资源内存（psutil.virtual\_memory.total）进行计算，与 lookahead cache 的 `--cpu-decoding-memory-utilization`（默认0.15）类似。请预留足够内存或配置相应比例值。
- `--apc-offload-not-lazy`
  - 是否关闭 offload 按层延迟启动，仅作调试用途。
- `--apc-offload-min-access-threshold`
  - 一个 block 会被 offload 的最小使用阈值，默认为1，即所有 block 都会被 offload。增加此值，被多次使用的 block 才会被 offload。
- `--apc-offload-enable-hit-cnt`
  - 打开 prefix cache offload 的命中率 log，每100个 block 打印一次。
- `--apc-offload-gpu-evictor-limit`
  - gpu evictor 的 free table 大小，默认为-1不生效，设置具体值限制 gpu上prefix cache 容量，仅作调试用途。

## 场景

- Auto Prefix Caching 适用于 common prefix 较多，prefill 计算占比较大的场景，例如 多轮对话，system prompt，代码补全等等。此类场景打开 `--enable-prefix-caching` 即可加速命中 gpu prefix cache 的请求的 prefill 阶段。
- GPU 容纳 prefix cache 的空间是有限的，当 common prefix 的累积量比较多，相同 common prefix 的请求间相隔比较远的场景下，之前保留的 prefix cache 已经被驱逐而无法获得收益。
  - 比如 GPU 的 block 数为100，有11个不同 prefix 的请求[Q1, ..., Q11]，每个对应 prefix block 数为 10，则Q11结束后Q1的 prefix cache 会被驱逐。即使Q12跟Q1的 prefix 一致，也无法获得加速。
  - 此时 prefix cache offload 通过额外的内存空间保留被驱逐的 prefix cache，进而加速Q12。
  - 可开启 `--apc-offload-enable-hit-cnt` 参数，通过 log cpu 判断是否有 offload 的 prefix cache 得到命中。该 log 统计所有 allocate的block 的 hit 情况。

[HIT CNT] total: 177800, gpu: 21944 (12.34%) cpu: 66166 (37.21%) not hit: 89690 (50.44%)

# 量化

在本文档中，我们将介绍模型量化的基本概念，以及使用 TACO-LLM 部署量化模型的完整流程实践。

## 量化概述

模型量化通常是指将一个连续取值(通常是 fp32, fp16)或者大量离散值的浮点型权重，转化为有限个离散值(通常是 int8, int4)的过程。这个过程会带来轻微的推理损失精度，但是存在如下优势：

- 减小模型体积
- 降低内存占用
- 在支持低精度运算的设备上提升推理速度

### 1. 量化比特

工业界目前常用的量化比特位数是 4 bits 和 8 bits，低于 4bits 的量化位宽被称为低比特量化。

### 2. 量化目标

- 权重：权重的量化是最常规的，量化权重可以减少模型大小和占用空间。
- 激活：量化激活可以大大减少内存占用，结合权重的量化可以充分利用设备的算力。
- KV cache：显存占用会随着生成的序列长度线性增长，量化 KV cache 可以节省显存，从而能够处理更大批次的大小。  
量化还可以选择不同的量化粒度，例如 per-tensor, per-group 等等。并且对于激活还有动态量化和静态量化的区别。

### 3. 量化形式

- 线性量化：将浮点数值域均匀的映射到整数值域，用固定的步长进行量化。该方式实现简单，硬件比较友好，适合分布相对均匀的数据。
- 非线性量化：根据数据的实际分布特征进行不均匀的量化，在数据密集区域使用更细的量化粒度。实现和计算都比较复杂，理论上可以获得更好的量化精度。  
在实际的推理业务中，由于非线性量化的计算复杂度较高，通常使用线性量化的方式。

### 4. 量化方法

- 量化感知训练 (Quantization Aware Training, QAT)：在训练过程中模拟量化效果，通过反向传播来补偿量化误差，让模型适应量化带来的损失。
- 训练后量化 (Post Training Quantization, PTQ)：在模型训练完成后，使用少量校准数据来确定量化参数，直接将模型量化，无需重新训练。  
在实际推理业务中，PTQ 的应用更加广泛。PTQ 的主要优势在于简单和高效，但可能会引入一定程度的精度损

失。

## TACO-LLM 量化支持

下面展示了 TACO-LLM 在各种硬件上对不同量化方案的支持情况：

- GPTQ：在 Volta, Turing, Ampere, Ada, Hopper, Intel CPU 上支持。
- AWQ：在 Turing, Ampere, Ada, Hopper, Intel CPU 上支持。
- Marlin：在 Ampere, Ada, Hopper 上支持。
- FP8：在 Ada, Hopper 上支持。
- Bitsandbytes：在 Turing, Ampere, Ada, Hopper 上支持。
- INT8(W8A8)：在 Turing, Ampere, Ada, Hopper 上支持。
- AQLM：在 Volta, Turing, Ampere, Ada, Hopper 上支持。

## TACO-LLM 快速启动

执行 `taco_llm serve -h` 命令可以查看 `taco-llm` 完整的在线模式参数配置，其中找到 `quantization` 的配置参数如下：

```
--quantization {aqlm,awq,deepspeedfp,tpu_int8,fp8,fbgemm_fp8,modelopt,marlin,gguf,gptq_marlin_24,gptq_marlin,awq_marlin,gptq,compressed-tensors,bitsandbytes,experts_int8,qqq,neuron_quant,None}, -q {aqlm,awq,deepspeedfp,tpu_int8,fp8,fbgemm_fp8,modelopt,marlin,gguf,gptq_marlin_24,gptq_marlin,awq_marlin,gptq,compressed-tensors,bitsandbytes,experts_int8,qqq,neuron_quant,None}
Method used to quantize the weights. If None, we first check the `quantization_config` attribute in the model config file. If that is None, we assume the model weights are not quantized and use `dtype` to determine the data type of the weights.
```

### 1. GPTQ-Marlin (AWQ)：

首先使用 AutoGPTQ(AutoAWQ) 将 fp16 模型权重量化。启动 TACO-LLM 的时候无需传入其他参数，server 会自动读取 config 文件中的量化参数来加载模型。TACO-LLM 在条件允许的情况下会默认使用 marlin kernel，可以传入 `--quantization gptq` 参数来强制使用 gptq kernel。

### 2. Bitsandbytes：

启动时添加启动参数 `--quantization bitsandbytes`，server 会自动读取 config 文件中的量化参数来加载模型。

### 3. FP8 (W8A8)：

TACO-LLM 采用动态量化的方案来将 BF16/FP16 量化到 FP8，并且不需要额外的矫正数据集。除了 `lm_head` 的所有 linear modules 都会按照 per-tensor 的方式进行量化。

```
from taco_llm import LLM
model = LLM(moth_path, quantization="fp8")
result = model.generate("Tell me about computer science.")
```

## GPTQ 量化实践 ( W4A16 ) :

本节以 TinyLlama-1.1B-Chat-v1.0 量化流程为例，介绍整个量化过程。

### 模型量化流程

1. 首先安装 autogptq，来作为量化工具。然后下载对应的模型权重 [TinyLlama-1.1B](#)。

```
pip install autogptq datasets transformers
```

2. 接下来可以使用下面脚本，来执行整个量化过程（其中所需的矫正数据集会自动下载）。矫正数据集可以优先使用模型对应的垂类数据集，如果没有的话，可以使用模型的预训练数据集或者是微调数据集。

```
import torch
from datasets import load_dataset
from gptqmodel import GPTQModel, QuantizeConfig
from transformers import AutoTokenizer

pretrained_model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
quantized_model_id = "TinyLlama-1.1B-Chat-v1.0-4bit-128g"

# os.makedirs(quantized_model_dir, exist_ok=True)
def get_wikitext2(tokenizer, nsamples, seqlen):
    traindata = load_dataset("wikitext", "wikitext-2-raw-v1", split="train").filter(
        lambda x: len(x["text"]) >= seqlen)

    return [tokenizer(example["text"]) for example in traindata.select(range(nsamples))]

@torch.no_grad()
def calculate_avg_ppl(model, tokenizer):
    from gptqmodel.utils import Perplexity

    ppl = Perplexity(
        model=model,
        tokenizer=tokenizer,
        dataset_path="wikitext",
        dataset_name="wikitext-2-raw-v1",
```

```
        split="train",
        text_column="text",
    )

    all = ppl.calculate(n_ctx=512, n_batch=512)

    # average ppl
    avg = sum(all) / len(all)

    return avg

def main():
    tokenizer = AutoTokenizer.from_pretrained(pretrained_model_id, use_fast=True)

    traindataset = get_wikitext2(tokenizer, nsamples=256, seqlen=1024)

    quantize_config = QuantizeConfig(
        bits=4, # quantize model to 4-bit
        group_size=128, # it is recommended to set the value to 128
        desc_act= False, #
    )

    # load un-quantized model, the model will always be force loaded into cpu
    model = GPTQModel.from_pretrained(pretrained_model_id, quantize_config)

    # quantize model, the calibration_dataset should be list of dict whose keys can only be "input_
ids" and "attention_mask"
    # with value under torch.LongTensor type.
    model.quantize(traindataset)

    # save quantized model
    model.save_quantized(quantized_model_id)

    # save quantized model using safetensors
    model.save_quantized(quantized_model_id, use_safetensors=True)

    # load quantized model, currently only support cpu or single gpu
    model = GPTQModel.from_quantized(quantized_model_id, device="cuda:0")

    # inference with model.generate
    print(tokenizer.decode(model.generate(**tokenizer("test is", return_tensors="pt").to("cuda:0"))
[0])))

    print(f"Quantized Model {quantized_model_id} avg PPL is {calculate_avg_ppl(model, tokeniz
r)}")
```

```
if __name__ == "__main__":
    import logging

    logging.basicConfig(
        format="%(asctime)s %(levelname)s [%(name)s] %(message)s",
        level=logging.INFO,
        datefmt="%Y-%m-%d %H:%M:%S",
    )

    main()
```

下面介绍一下量化参数的选择：

- `--bits`：权重量化的位宽。根据需求选择，要求节省显存选择 4，但是会对精度有较大的影响。基于显存和精度的平衡，建议选择 8，此时精度基本没有损失。
- `--group_size`：group 量化的 size 大小，越小精度越高，但是会增加推理成本。建议选择 128。
- `--desc_act`：是否使用激活重排。打开会提高量化精度，但是会增加推理成本。建议选择 False。
- `--nsamples`：矫正数据集的样本数量。数量太多会增加量化时间，且还会改变权重分布。建议选择 256。
- `--seqlen`：矫正数据集的样本长度。数量太多会增加量化时间，且还会改变权重分布。7B 模型建议选择 2048，70B 以上模型建议选择 4096。

```
bits = [4,8]
group_size = [64,128]
nsamples = [256,512]
seqlen = [2048, 4096]
desc_act = [True, False]
```

# CPU 辅助加速

## 介绍

大模型的自回归解码的特性导致其不能充分利用 GPU 的并行计算进行加速。基于此，学术界提出了投机采样，其可以利用更高效的 Draft Model 快速生成多个 Token，然后一次性地交给 Target Model 进行验证。在接受率高的情况下，可以显著减少推理时间。

传统的投机采样使用 GPU 作为 Draft Model 的计算资源，而 GPU 的成本高昂。针对这个痛点，TACO 与 Intel 团队合作，基于 AMX 指令集对 CPU 上的矩阵乘法做了优化，借助 Intel 推出的 IPEX 加速库及其他技术对 CPU 上的推理进行加速，使得 CPU 作为 Draft Model 成为可能，从而在进行推理加速的同时，显著降低了推理成本。

注意：

当前该功能支持的机型：搭载 Intel EMR 8576C 或 Intel SPR 8476C 的 GPU 云服务器，其中：

- SPR：搭载内置加速器的第四代英特尔® 至强® 可扩展处理器。
- EMR：搭载内置加速器的第五代英特尔® 至强® 可扩展处理器。

## 基于 TACO-LLM 使用 CPU 单独进行推理

注意：

该 feature 会导致 GPU 推理的功能失效，因此若想使用 GPU 推理，需要重新安装对应的包。

### 1. 额外依赖安装

```
sudo apt-get update
sudo apt-get install -y libdnnl-dev
pip install intel-extension-for-pytorch==2.4.0
pip install torch --index-url https://download.pytorch.org/whl/cpu
```

### 2. 环境检查 (如果出现问题)

需要确保：

1. torch 和 ipex 的版本号匹配
2. torch 是 cpu 版本

```
root@c2d5f7048105:/script# pip show torch
Name: torch
Version: 2.4.0+cpu
Summary: Tensors and Dynamic neural networks in Python with strong GPU acceleration
Home-page: https://pytorch.org/
Author: PyTorch Team
Author-email: packages@pytorch.org
License: BSD-3
Location: /usr/local/lib/python3.10/dist-packages
Requires: filelock, fsspec, Jinja2, networkx, sympy, typing-extensions
Required-by: compressed-tensors, flash-attn, lightning-thunder, tensorizer, torch-tensorrt, torchaudio, torchdata, torchtext, torchvision, vllm, vllm-flash-attn, xformers
root@c2d5f7048105:/script# pip show intel-extension-for-pytorch
Name: intel_extension_for_pytorch
Version: 2.4.0
Summary: Intel® Extension for PyTorch*
Home-page: https://github.com/intel/intel-extension-for-pytorch
Author: Intel Corp.
Author-email:
License: https://www.apache.org/licenses/LICENSE-2.0
Location: /usr/local/lib/python3.10/dist-packages
Requires: numpy, packaging, psutil
Required-by:
root@c2d5f7048105:/script#
```

### 3. 使用

以离线模式为例，其使用方法和在 GPU 上推理完全一致。

```
llm = LLM(model="facebook/opt-125m")
```

## 基于 TACO-LLM 使用 CPU 辅助投机采样

### 1. 额外依赖安装

```
sudo apt-get update
sudo apt-get install -y libdnnl-dev
pip install intel-extension-for-pytorch==2.4.0
pip install torch==2.4.0
```

### 2. 环境检查 (如果出现问题)

需要确保：

1. torch 和 ipex 的版本号匹配
2. torch 是不带 cpu 的版本

```

● root@b24522df0352:~# pip show torch
Name: torch
Version: 2.4.0
Summary: Tensors and Dynamic neural networks in Python with strong GPU acceleration
Home-page: https://pytorch.org/
Author: PyTorch Team
Author-email: packages@pytorch.org
License: BSD-3
Location: /usr/local/lib/python3.10/dist-packages
Requires: filelock, fsspec, Jinja2, networkx, nvidia-cublas-cu12, nvidia-cuda-cupti-cu12, nvidia-cuda-nvrtc-cu12, nvidia-cuda-runtime-cu12, nvidia-cudnn-cu12, nvidia-cufft-cu12, nvidia-curand-cu12, nvidia-cusolver-cu12, nvidia-cuspars-cu12, nvidia-nccl-cu12, nvidia-nvtx-cu12, sympy, triton, typing-extensions
Required-by: accelerate, compressed-tensors, flash-attn, lightning-thunder, peft, sentence-transformers, taco-llm, tensorizer, timm, torch-tensorrt, torchdata, torchtext, torchvision, vllm-flash-attn, xformers
● root@b24522df0352:~# pip show intel-extension-for-pytorch
Name: intel_extension_for_pytorch
Version: 2.4.0
Summary: Intel® Extension for PyTorch*
Home-page: https://github.com/intel/intel-extension-for-pytorch
Author: Intel Corp.
Author-email:
License: https://www.apache.org/licenses/LICENSE-2.0
Location: /usr/local/lib/python3.10/dist-packages
Requires: numpy, packaging, psutil
Required-by:

```

### 3. 使用

以离线模式为例，其使用方法在投机采样原始配置上额外新增一个 `cpu_draft_worker` 即可。

```

llm = LLM(
    model = "meta-llama/Llama-2-7b-chat-hf",
    speculative_model = "Felladrin/Llama-68M-Chat-v1",
    num_speculative_tokens = 2,
    use_v2_block_manager = True,
    cpu_draft_worker = True,          # < < == 新增参数
)

```

## 附录

### AMX 介绍

#### 什么是 AMX ?

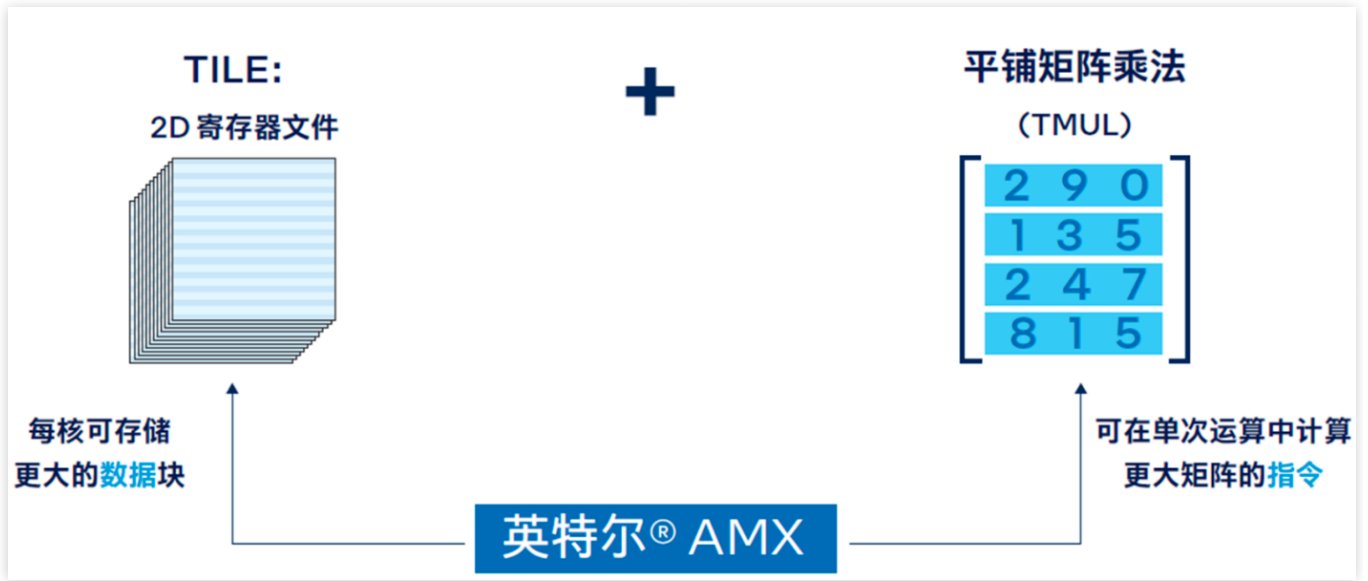
英特尔推出的第四代英特尔® 至强® 可扩展处理器及其内置的英特尔® 高级矩阵扩展 (Intel® Advanced Matrix Extensions, 英特尔® AMX) 可进一步提高 AI 功能，实现较上一代产品 3 至 10 倍的推理和训练性能提升。

开发人员可以编写非 AI 功能代码来利用处理器的指令集架构 (ISA)，也可编写 AI 功能代码，以充分发挥英特尔® AMX 指令集的优势。英特尔已将其 oneAPI DL 引擎——英特尔® oneAPI 深度神经网络库 (Intel® oneAPI Deep Neural Network Library, 英特尔® oneDNN) 集成至包括 TensorFlow、PyTorch、PaddlePaddle 和 ONNX 在内的多个主流 AI 应用开源工具当中。

#### AMX 架构

英特尔® AMX 架构由两部分组件构成：

- 第一部分为 TILE，由 8 个 1 KB 大小的 2D 寄存器组成，可存储大数据块。
- 第二部分为平铺矩阵乘法 (TMUL)，它是与 TILE 连接的加速引擎，可执行用于 AI 的矩阵乘法计算。



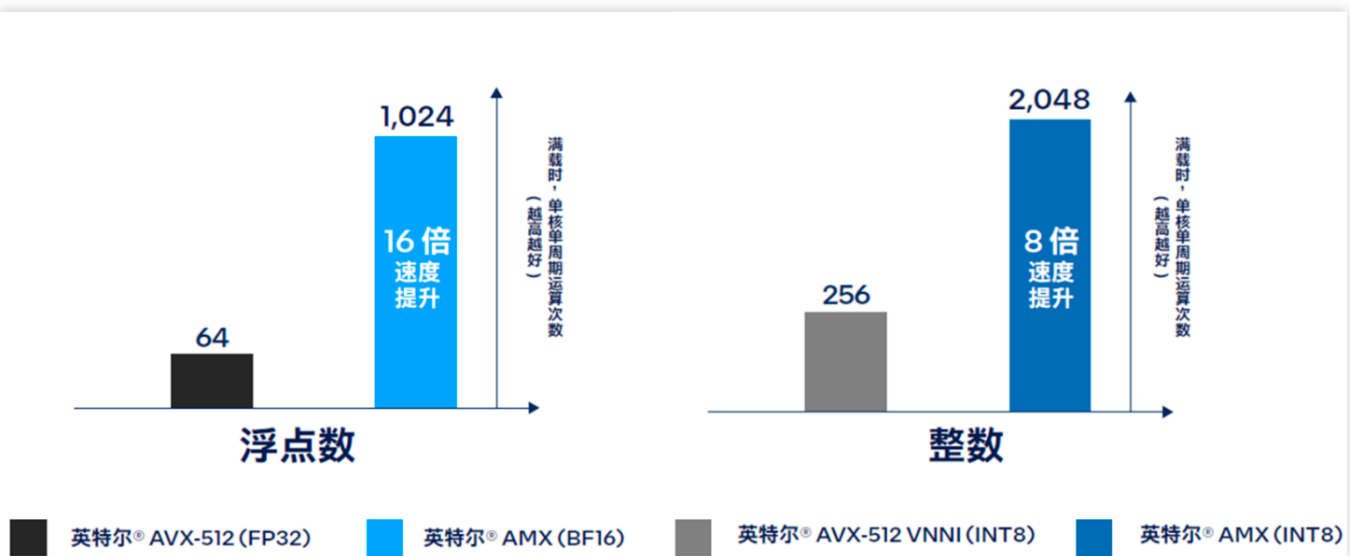
### AMX 支持的数据类型

英特尔® AMX 支持两种数据类型：INT8 和 BF16，两者均可用于 AI 工作负载所需的矩阵乘法运算。

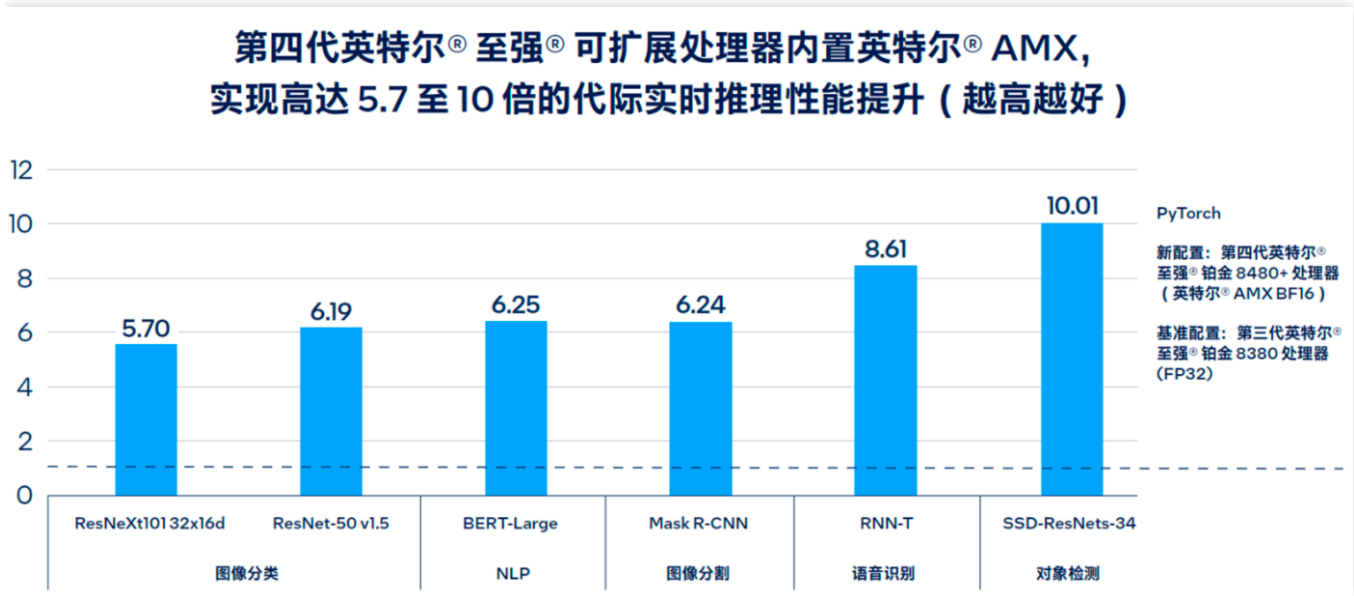
- 当推理无需 FP32 的精度时可使用 INT8 这种数据类型。由于该数据类型的精度较低，因此单位计算周期内运算次数就更多。
- BF16 这种数据类型实现的准确度足以达到大多数训练的要求，必要时它也能让 AI 推理实现更高的准确度。

### AMX 的性能

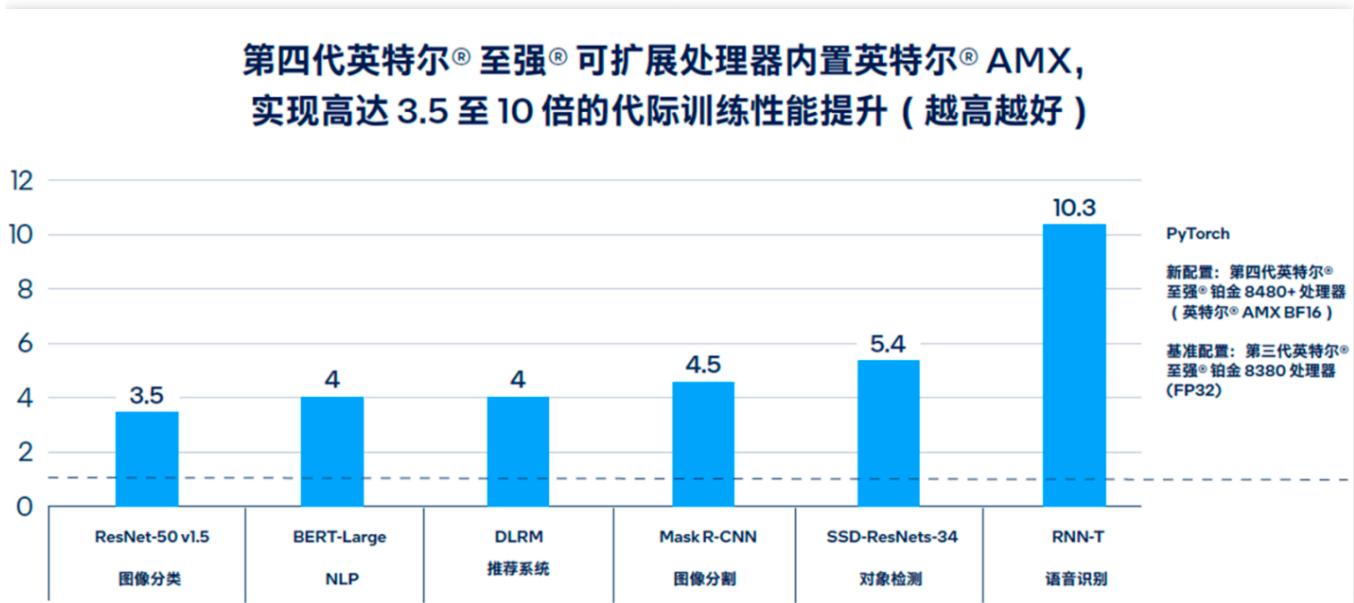
凭借这种新的平铺架构，英特尔® AMX 实现了大幅代际性能提升。与运行英特尔® 高级矢量扩展 512 神经网络指令 (Intel® Advanced Vector Extensions 512 Neural Network Instructions, 英特尔® AVX-512 VNNI) 的第三代英特尔® 至强® 可扩展处理器相比，运行英特尔® AMX 的第四代英特尔® 至强® 可扩展处理器将单位计算周期内执行 INT8 运算的次数从 256 次提高至 2048 次。此外，如图 6 所示，第四代英特尔® 至强® 可扩展处理器可在单位计算周期内执行 1024 次 BF16 运算，而第三代英特尔® 至强® 可扩展处理器执行 FP32 运算的次数仅为 64 次。



下图所示为英特尔® AMX 在代际间实现高达 5.7 至 10 倍的 PyTorch 实时推理性能提升的情况。



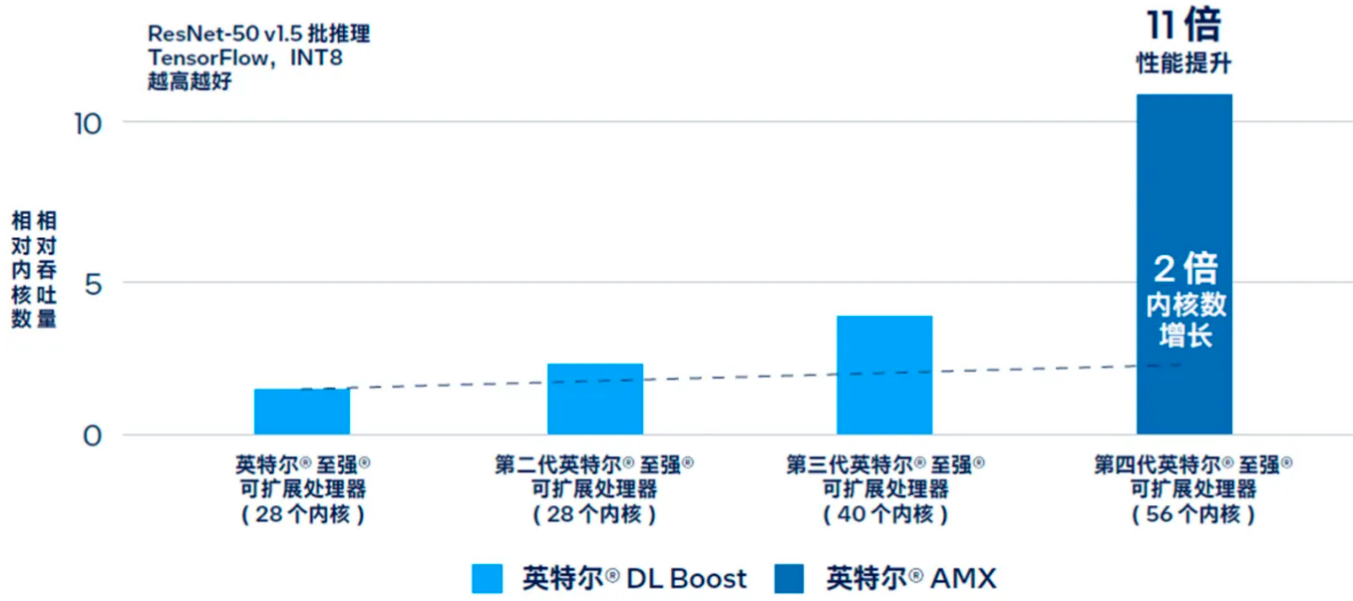
下图所示为英特尔® AMX 在代际间实现高达 3.5 至 10 倍的 PyTorch 训练性能提升的情况。



通过下图可以看出英特尔® AMX 带来的性能提升远大于每一代产品 ( 从第一代英特尔® 至强® 可扩展处理器开始 ) 通过增加内核所实现的性能提升。

## 摩尔定律与加速器

为工作负载匹配合适的计算引擎



### IPEX 介绍

Intel® Extension for PyTorch\* (IPEX) 是由英特尔发起的一个开源扩展项目，旨在通过模块级的全面优化及更为简洁的API接口，在基于原生PyTorch框架运行时，显著提升深度学习任务在英特尔硬件（包括但不限于CPU和GPU）上的推理与训练性能。IPEX兼容PyTorch生态系统中超过90%的主流模型，并对其中50多个深度模型进行了特别优化。用户只需添加少量代码以启用BF16混合精度支持，即可轻松享受到性能上的显著改进，整个过程无需复杂的配置调整，从而提供了近乎即插即用的便捷体验。

此外，Intel® Extension for PyTorch\* 通过对Intel硬件特性的深入挖掘，如利用Intel® Advanced Vector Extensions 512 (AVX-512) 中的向量神经网络指令(VNNI)、Intel® Advanced Matrix Extensions (AMX) 以及Intel独立显卡上配备的Intel Xe Matrix Extensions (XMX) AI加速引擎等技术，进一步增强了PyTorch在Intel平台上的执行效率。

# 长序列优化

## 背景

在 LLM 大模型推理中，长序列场景应用越来越广泛，目前业界对长序列的优化主要是 kv cache 量化、稀疏化等方法，这些都对模型精度有一定的影响。TACO LLM 的长序列并行方案，可以在长序列场景进行精度无损的加速。长序列推理场景首字延迟会比较慢，针对该问题，序列并行在 prefill 阶段采用了 Ring Attention 类的方案，可以通过扩展机器，降低首字延迟。而对于推理阶段，TACO LLM 可以使用 Lookahead 等投机采样技术加速。

## 使用方式

启动 TACO LLM 服务时添加参数 `--sequence-parallel-size 2` 即可开启序列并行，目前只支持 `serve` 方式，不支持 LLM 的离线方式（0.6.4版本开始支持）。

## 最佳实践

序列并行可以和张量并行一起使用，同时也可以和 FP8 一起使用。序列并行适合符合如下几个特点的场景：

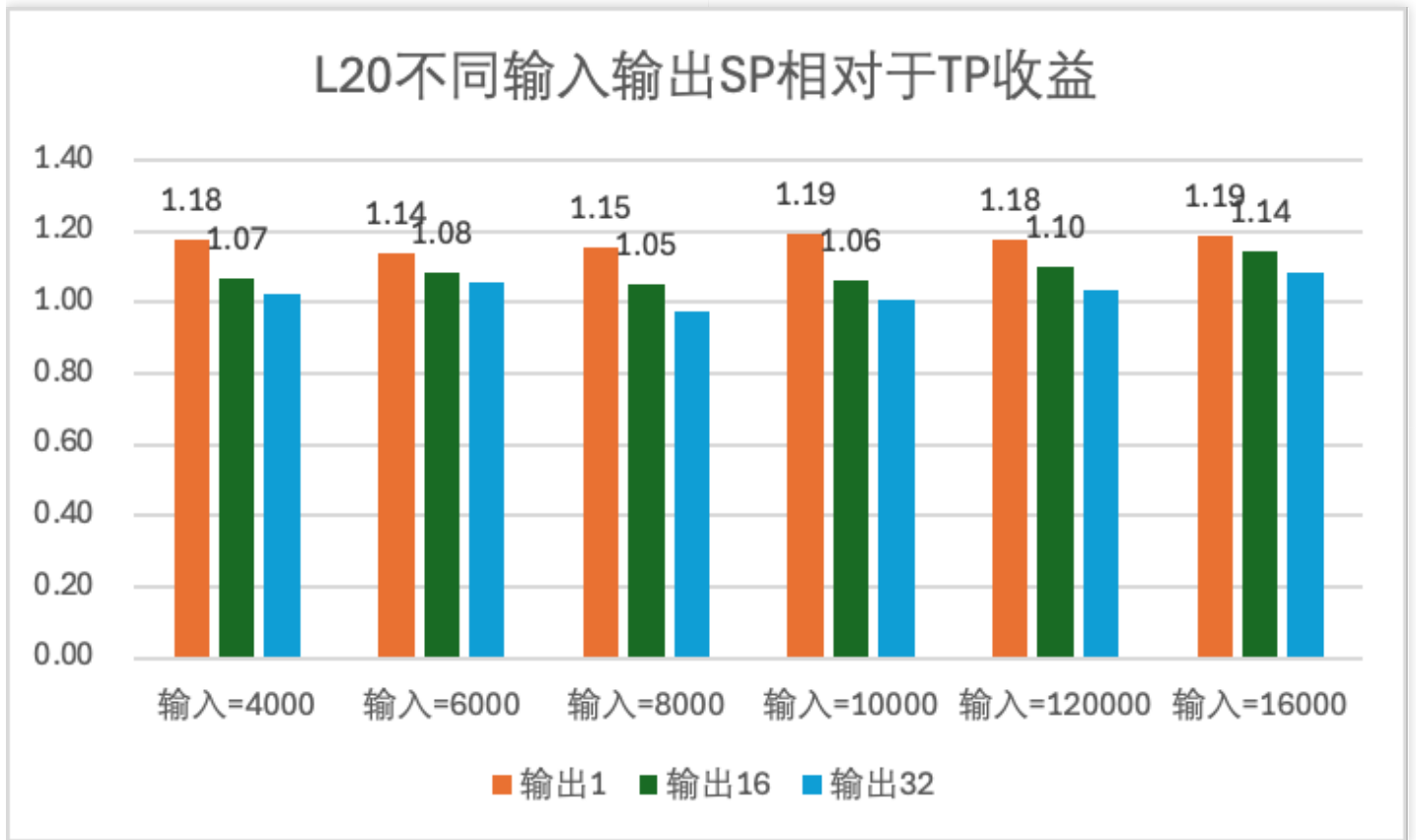
- 由于序列并行具备通信计算重叠的优点，故适合在 L20、4090 等 PCIE 系列卡上。
- 适合输入较长，输出较短的场景。
- 建议与 FP8 一起使用，加速效果更好。
- 适合 GQA 类的模型，由于 kv cache 更小，通信更小，相对于 TP 来说加速更快。

## 性能数据

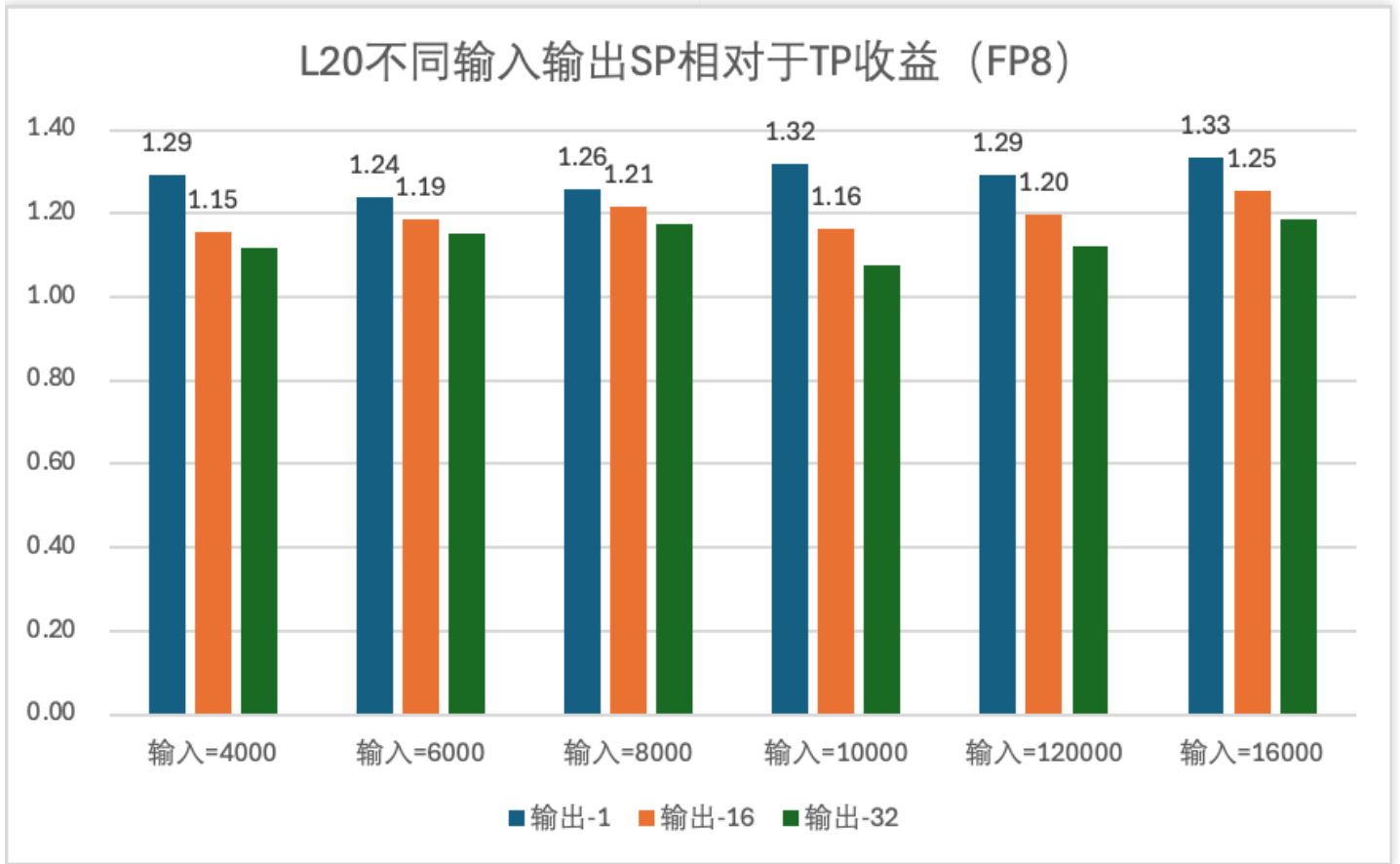
### MHA 场景

当前测试模型大小为 Llama2 7B，序列长度为6000，SP 为2，在H20和A800上，相对于 TP 无性能提升。

在 L20上FP16：SP = 2的场景相对于 TP = 2场景，端到端性能有5%到15%的提升，prefill 性能提升15%-20%。



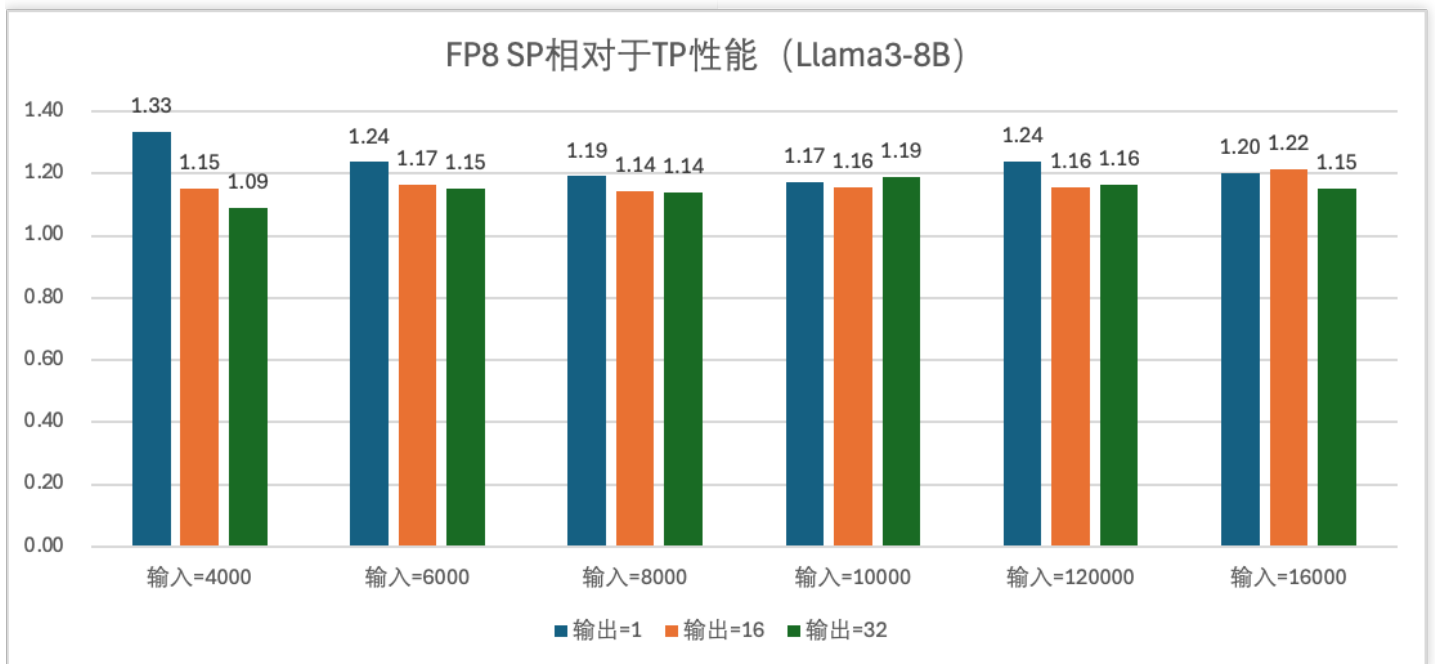
在 L20上 FP8 : SP = 2的场景相对于 TP = 2场景，性能有15%到25%的提升。



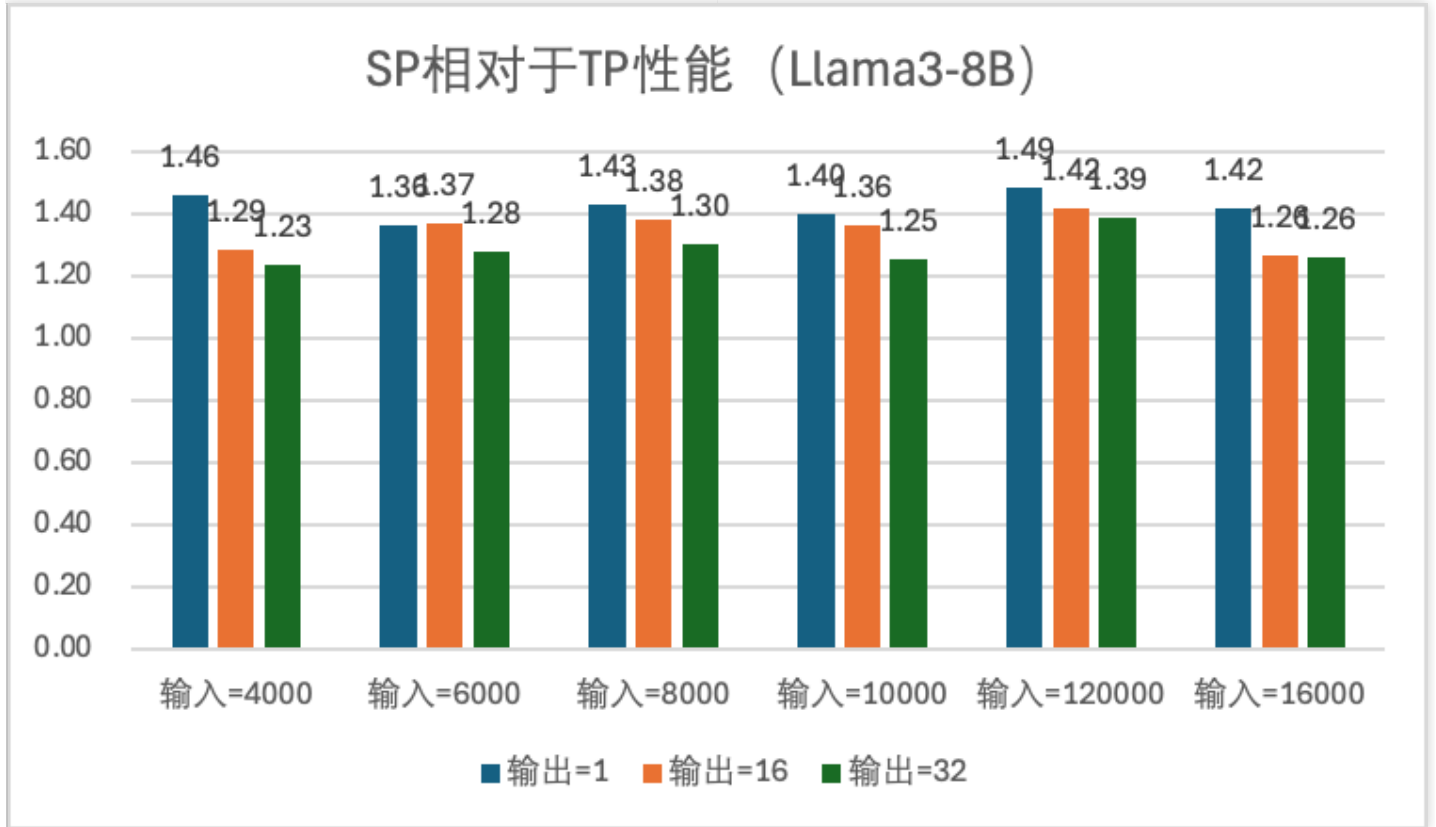
## GQA 场景

当前测试模型大小为 Llama3-8B。

在 L20上FP16 : SP = 2的场景相对于TP = 2场景，端到端性能有10%到20%的提升，prefill 性能提升20%-30%。



在 L20上FP8 : SP = 2的场景相对于 TP = 2场景，端到端性能有20%到40%的提升，prefill 性能有40%-50%。



# TACO LLM API

## Offline API

### LLM

#### LLM 构造参数

```
class taco_llm.LLM(  
    model: str,  
    tokenizer: Optional[str] = None,  
    tokenizer_mode: str = "auto",  
    skip_tokenizer_init: bool = False,  
    trust_remote_code: bool = False,  
    tensor_parallel_size: int = 1,  
    dtype: str = "auto",  
    quantization: Optional[str] = None,  
    revision: Optional[str] = None,  
    tokenizer_revision: Optional[str] = None,  
    seed: int = 0,  
    gpu_memory_utilization: float = 0.9,  
    swap_space: float = 4,  
    cpu_offload_gb: float = 0,  
    enforce_eager: Optional[bool] = None,  
    max_context_len_to_capture: Optional[int] = None,  
    max_seq_len_to_capture: int = 8192,  
    disable_custom_all_reduce: bool = False,  
    disable_async_output_proc: bool = False,  
    **kwargs,  
)  
"""
```

This class includes a tokenizer, a language model (possibly distributed across multiple GPUs), and GPU memory space allocated for intermediate states (aka KV cache). Given a batch of prompts and sampling parameters, this class generates texts from the model, using an intelligent batching mechanism and efficient memory management.

#### Args:

**model:** The name or path of a HuggingFace Transformers model.

**tokenizer:** The name or path of a HuggingFace Transformers tokenizer.

**tokenizer\_mode:** The tokenizer mode. "auto" will use the fast tokenizer if available, and "slow" will always use the slow tokenizer.

**skip\_tokenizer\_init:** If true, skip initialization of tokenizer and detokenizer. Expect valid prompt\_token\_ids and None for prompt

- from the input.
- `trust_remote_code`: Trust remote code (e.g., from HuggingFace) when downloading the model and tokenizer.
- `tensor_parallel_size`: The number of GPUs to use for distributed execution with tensor parallelism.
- `dtype`: The data type for the model weights and activations. Currently, we support float32, float16, and bfloat16. If auto, we use the `torch_dtype` attribute specified in the model config file. However, if the `torch_dtype` in the config is float32, we will use float16 instead.
- `quantization`: The method used to quantize the model weights. Currently, we support "awq", "gptq", and "fp8" (experimental). If None, we first check the `quantization_config` attribute in the model config file. If that is None, we assume the model weights are not quantized and use `dtype` to determine the data type of the weights.
- `revision`: The specific model version to use. It can be a branch name, a tag name, or a commit id.
- `tokenizer_revision`: The specific tokenizer version to use. It can be a branch name, a tag name, or a commit id.
- `seed`: The seed to initialize the random number generator for sampling.
- `gpu_memory_utilization`: The ratio (between 0 and 1) of GPU memory to reserve for the model weights, activations, and KV cache. Higher values will increase the KV cache size and thus improve the model's throughput. However, if the value is too high, it may cause out-of-memory (OOM) errors.
- `swap_space`: The size (GiB) of CPU memory per GPU to use as swap space. This can be used for temporarily storing the states of the requests when their best\_of sampling parameters are larger than 1. If all requests will have `best_of=1`, you can safely set this to 0. Otherwise, too small values may cause out-of-memory (OOM) errors.
- `cpu_offload_gb`: The size (GiB) of CPU memory to use for offloading the model weights. This virtually increases the GPU memory space you can use to hold the model weights, at the cost of CPU-GPU data transfer for every forward pass.
- `enforce_eager`: Whether to enforce eager execution. If True, we will disable CUDA graph and always execute the model in eager mode. If False, we will use CUDA graph and eager execution in hybrid.
- `max_context_len_to_capture`: Maximum context len covered by CUDA graphs. When a sequence has context length larger than this, we fall back to eager mode (DEPRECATED. Use `max_seq_len_to_capture` instead).
- `max_seq_len_to_capture`: Maximum sequence len covered by CUDA graphs. When a sequence has context length larger than this, we fall back to eager mode.
- `disable_custom_all_reduce`: See `ParallelConfig`
- \*\*kwargs: Arguments for `:class:taco_llm.EngineArgs`.

**Note:**

This class is intended to be used for offline inference. For online serving, use the `:class:taco_llm.AsyncLLMEngine` class instead.

```
"""
```

TACO-LLM 支持离线和在线两种模式，这两种模式的参数配置是一致的。因此，除了上述明确提到的参数外，您还可以设置任意 TACO-LLM 在线模式支持的参数。完整的参数配置请参见 [Online API 章节](#)。

**chat 接口**

```
def chat(
    self,
    messages: List[ChatCompletionMessageParam],
    sampling_params: Optional[Union[SamplingParams,
                                   List[SamplingParams]]] = None,
    use_tqdm: bool = True,
    lora_request: Optional[LoRARequest] = None,
    chat_template: Optional[str] = None,
    add_generation_prompt: bool = True,
) -> List[RequestOutput]:
    """
```

Generate responses for a chat conversation.

The chat conversation is converted into a text prompt using the tokenizer and calls the `:meth:`generate`` method to generate the responses.

Multi-modal inputs can be passed in the same way you would pass them to the OpenAI API.

**Args:**

`messages`: A single conversation represented as a list of messages.

Each message is a dictionary with 'role' and 'content' keys.

`sampling_params`: The sampling parameters for text generation.

If None, we use the default sampling parameters. When it is a single value, it is applied to every prompt. When it is a list, the list must have the same length as the prompts and it is paired one by one with the prompt.

`use_tqdm`: Whether to use tqdm to display the progress bar.

`lora_request`: LoRA request to use for generation, if any.

`chat_template`: The template to use for structuring the chat.

If not provided, the model's default chat template will be used.

`add_generation_prompt`: If True, adds a generation template to each message.

**Returns:**

A list of ``RequestOutput`` objects containing the generated responses in the same order as the input messages.

"""

## generate 接口

```
def generate(
    self,
    prompts: Union[Union[PromptInputs, Sequence[PromptInputs]],
                  Optional[Union[str, List[str]]]] = None,
    sampling_params: Optional[Union[SamplingParams,
                                   Sequence[SamplingParams]]] = None,
    prompt_token_ids: Optional[Union[List[int], List[List[int]]]] = None,
    use_tqdm: bool = True,
    lora_request: Optional[Union[List[LoRARequest], LoRARequest]] = None,
    prompt_adapter_request: Optional[PromptAdapterRequest] = None,
    guided_options_request: Optional[Union[LLMGuidedOptions,
                                           GuidedDecodingRequest]] = None
) -> List[RequestOutput]:
```

"""Generates the completions for the input prompts.

This class automatically batches the given prompts, considering the memory constraint. For the best performance, put all of your prompts into a single list and pass it to this method.

### Args:

inputs: A list of inputs to generate completions for.

sampling\_params: The sampling parameters for text generation. If None, we use the default sampling parameters.

When it is a single value, it is applied to every prompt.

When it is a list, the list must have the same length as the prompts and it is paired one by one with the prompt.

use\_tqdm: Whether to use tqdm to display the progress bar.

lora\_request: LoRA request to use for generation, if any.

prompt\_adapter\_request: Prompt Adapter request to use for generation, if any.

### Returns:

A list of ``RequestOutput`` objects containing the generated completions in the same order as the input prompts.

### Note:

Using ``prompts`` and ``prompt\_token\_ids`` as keyword parameters is considered legacy and may be deprecated in the future. You should instead pass them via the ``inputs`` parameter.

"""

# Online API

执行命令 `taco_llm serve -h` 可以查看 TACO-LLM 的完整在线模式参数配置：

```
# taco_llm serve -h
```

```
usage: taco_llm serve <model_tag> [options]
```

positional arguments:

```
model_tag          The model tag to serve
```

options:

```
-h, --help          show this help message and exit
--config CONFIG     Read CLI options from a config file. Must be a YAML with the following options:
--host HOST         host name
--port PORT         port number
--uvicorn-log-level {debug,info,warning,error,critical,trace}
                   log level for uvicorn
--allow-credentials allow credentials
--allowed-origins ALLOWED_ORIGINS
                   allowed origins
--allowed-methods ALLOWED_METHODS
                   allowed methods
--allowed-headers ALLOWED_HEADERS
                   allowed headers
--api-key API_KEY  If provided, the server will require this key to be presented in the header.
--lora-modules LORA_MODULES [LORA_MODULES ...]
                   LoRA module configurations in the format name=path. Multiple modules can be specified.
--prompt-adapters PROMPT_ADAPTERS [PROMPT_ADAPTERS ...]
                   Prompt adapter configurations in the format name=path. Multiple adapters can be specified.
--chat-template CHAT_TEMPLATE
                   The file path to the chat template, or the template in single-line form for the specified model
--response-role RESPONSE_ROLE
                   The role name to return if request.add_generation_prompt=true.
--ssl-keyfile SSL_KEYFILE
                   The file path to the SSL key file
--ssl-certfile SSL_CERTFILE
                   The file path to the SSL cert file
--ssl-ca-certs SSL_CA_CERTS
                   The CA certificates file
--ssl-cert-reqs SSL_CERT_REQS
                   Whether client certificate is required (see stdlib ssl module's
```

`--root-path ROOT_PATH`  
FastAPI `root_path` when app is behind a path based routing proxy

`--middleware MIDDLEWARE`  
Additional ASGI middleware to apply to the app. We accept multiple `--middleware` arguments. The value should be an import path. If a function is provided, `taco_llm` will add it to the server using `@app.middleware('http')`.  
If a class is provided, `taco_llm` will add it to the server using `app.add_middleware()`.

`--return-tokens-as-token-ids`  
When `--max-logprobs` is specified, represents single tokens as strings of the form `'token_id:{token_id}'` so that tokens that are not JSON-encodable can be identified.

`--disable-frontend-multiprocessing`  
If specified, will run the OpenAI frontend server in the same process as the model serving engine.

`--enable-auto-tool-choice`  
Enable auto tool choice for supported models. Use `--tool-call-parser` to specify which parser to use

`--tool-call-parser {mistral,hermes}`  
Select the tool call parser depending on the model that you're using. This is used to parse the model-generated tool call into OpenAI API format. Required for `--enable-auto-tool-choice`.

`--model MODEL` Name or path of the huggingface model to use.

`--tokenizer TOKENIZER`  
Name or path of the huggingface tokenizer to use. If unspecified, model name or path will be used.

`--skip-tokenizer-init`  
Skip initialization of tokenizer and detokenizer

`--revision REVISION` The specific model version to use. It can be a branch name, a tag name, or a commit id. If unspecified, will use the default version.

`--code-revision CODE_REVISION`  
The specific revision to use for the model code on Hugging Face Hub. It can be a branch name, a tag name, or a commit id. If unspecified, will use the default version.

`--tokenizer-revision TOKENIZER_REVISION`  
Revision of the huggingface tokenizer to use. It can be a branch name, a tag name, or a commit id. If unspecified, will use the default version.

`--tokenizer-mode {auto,slow,mistral}`  
The tokenizer mode. \* "auto" will use the fast tokenizer if available. \* "slow" will always use the slow tokenizer. \* "mistral" will always use the `mistral_common` tokenizer.

`--trust-remote-code` Trust remote code from huggingface.

`--download-dir DOWNLOAD_DIR`  
Directory to download and load the weights, default to the default cache dir of huggingface.

`--load-format {auto,pt,safetensors,np-cache,dummy,tokenizer,sharded_state,gguf,bitsandbytes,mistral}`  
The format of the model weights to load. \* "auto" will try to load the weights in the safetensors format and fall back to the pytorch bin format if safetensors format is not available. \* "pt" will load the weights in the  
pytorch bin format. \* "safetensors" will load the weights in the safetensors format. \* "np-cache" will load the weights in pytorch format and store a numpy cache to speed up the loading. \* "du

memory" will initialize the

weights with random values, which is mainly for profiling. \* "tensorizer" will load the weights using tensorizer from CoreWeave. See the Tensorize vLLM Model script in the Examples section for more information. \*

"bitsandbytes" will load the weights using bitsandbytes quantization.

--config-format {auto,hf,mistral}

The format of the model config to load. \* "auto" will try to load the config in hf format if available else it will try to load in mistral format

--dtype {auto,half,float16,bfloat16,float,float32}

Data type for model weights and activations. \* "auto" will use FP16 precision for FP32 and FP16 models, and BF16 precision for BF16 models. \* "half" for FP16. Recommended for AWQ quantization. \* "float16" is the same as

"half". \* "bfloat16" for a balance between precision and range. \* "float" is shorthand for FP32 precision. \* "float32" for FP32 precision.

--kv-cache-dtype {auto,fp8,fp8\_e5m2,fp8\_e4m3}

Data type for kv cache storage. If "auto", will use model data type. CUDA 11.8+ supports fp8 (=fp8\_e4m3) and fp8\_e5m2. ROCm (AMD GPU) supports fp8 (=fp8\_e4m3)

--quantization-param-path QUANTIZATION\_PARAM\_PATH

Path to the JSON file containing the KV cache scaling factors. This should generally be supplied, when KV cache dtype is FP8. Otherwise, KV cache scaling factors default to 1.0, which may cause accuracy issues. FP8\_E5M2

(without scaling) is only supported on cuda version greater than 11.8. On ROCm (AMD GPU), FP8\_E4M3 is instead supported for common inference criteria.

--max-model-len MAX\_MODEL\_LEN

Model context length. If unspecified, will be automatically derived from the model config.

--guided-decoding-backend {outlines,lm-format-enforcer}

Which engine will be used for guided decoding (JSON schema / regex etc) by default. Currently support outlines and lm-format-enforcer. Can be overridden per request via guided\_decoding\_backend parameter.

--distributed-executor-backend {ray,mp}

Backend to use for distributed serving. When more than 1 GPU is used, will be automatically set to "ray" if installed or "mp" (multiprocessing) otherwise.

--worker-use-ray Deprecated, use --distributed-executor-backend=ray.

--pipeline-parallel-size PIPELINE\_PARALLEL\_SIZE, -pp PIPELINE\_PARALLEL\_SIZE

Number of pipeline stages.

--tensor-parallel-size TENSOR\_PARALLEL\_SIZE, -tp TENSOR\_PARALLEL\_SIZE

Number of tensor parallel replicas.

--max-parallel-loading-workers MAX\_PARALLEL\_LOADING\_WORKERS

Load model sequentially in multiple batches, to avoid RAM OOM when using tensor parallel and large models.

--ray-workers-use-nsight

If specified, use nsight to profile Ray workers.

--block-size {8,16,32}

Token block size for contiguous chunks of tokens. This is ignored on neuron devices and set to max-model-len

--enable-prefix-caching

- Enables automatic prefix caching.
- `--disable-sliding-window`  
Disables sliding window, capping to sliding window size
  - `--use-v2-block-manager`  
Use BlockSpaceMangerV2.
  - `--num-lookahead-slots NUM_LOOKAHEAD_SLOTS`  
Experimental scheduling config necessary for speculative decoding. This will be replaced by speculative config in the future; it is present to enable correctness tests until then.
  - `--seed SEED` Random seed for operations.
  - `--swap-space SWAP_SPACE`  
CPU swap space size (GiB) per GPU.
  - `--cpu-offload-gb CPU_OFFLOAD_GB`  
The space in GiB to offload to CPU, per GPU. Default is 0, which means no offloading. Intuitively, this argument can be seen as a virtual way to increase the GPU memory size. For example, if you have one 24 GB GPU and set this to 10, virtually you can think of it as a 34 GB GPU. Then you can load a 13B model with BF16 weight, which requires at least 26GB GPU memory. Note that this requires fast CPU-GPU interconnect, as part of the model is loaded from CPU memory to GPU memory on the fly in each model forward pass.
  - `--gpu-memory-utilization GPU_MEMORY_UTILIZATION`  
The fraction of GPU memory to be used for the model executor, which can range from 0 to 1. For example, a value of 0.5 would imply 50% GPU memory utilization. If unspecified, will use the default value of 0.9.
  - `--num-gpu-blocks-override NUM_GPU_BLOCKS_OVERRIDE`  
If specified, ignore GPU profiling result and use this number of GPU blocks. Used for testing preemption.
  - `--max-num-batched-tokens MAX_NUM_BATCHED_TOKENS`  
Maximum number of batched tokens per iteration.
  - `--max-num-seqs MAX_NUM_SEQS`  
Maximum number of sequences per iteration.
  - `--max-logprobs MAX_LOGPROBS`  
Max number of log probs to return logprobs is specified in SamplingParams.
  - `--disable-log-stats` Disable logging statistics.
  - `--quantization {aqlm,awq,deepspeedfp,tpu_int8,fp8,fbgemm_fp8,modelopt,marlin,gguf,gptq_marlin_24,gptq_marlin,awq_marlin,gptq,compressed-tensors,bitsandbytes,experts_int8,qqq,neuron_quant,None}`, `-q {aqlm,awq,deepspeedfp,tpu_int8,fp8,fbgemm_fp8,modelopt,marlin,gguf,gptq_marlin_24,gptq_marlin,awq_marlin,gptq,compressed-tensors,bitsandbytes,experts_int8,qqq,neuron_quant,None}`  
Method used to quantize the weights. If None, we first check the `quantization_config` attribute in the model config file. If that is None, we assume the model weights are not quantized and use `dtype` to determine the data type of the weights.
  - `--rope-scaling ROPE_SCALING`  
RoPE scaling configuration in JSON format. For example, `{"type":"dynamic","factor":2.0}`
  - `--rope-theta ROPE_THETA`  
RoPE theta. Use with `rope_scaling`. In some cases, changing the RoPE theta improves the performance of the scaled model.
  - `--enforce-eager` Always use eager-mode PyTorch. If False, will use eager mode and CUDA graph in

hybrid for maximal performance and flexibility.

--max-context-len-to-capture MAX\_CONTEXT\_LEN\_TO\_CAPTURE

Maximum context length covered by CUDA graphs. When a sequence has context length larger than this, we fall back to eager mode. (DEPRECATED. Use --max-seq-len-to-capture instead)

--max-seq-len-to-capture MAX\_SEQ\_LEN\_TO\_CAPTURE

Maximum sequence length covered by CUDA graphs. When a sequence has context length larger than this, we fall back to eager mode.

--disable-custom-all-reduce

See ParallelConfig.

--tokenizer-pool-size TOKENIZER\_POOL\_SIZE

Size of tokenizer pool to use for asynchronous tokenization. If 0, will use synchronous tokenization.

--tokenizer-pool-type TOKENIZER\_POOL\_TYPE

Type of tokenizer pool to use for asynchronous tokenization. Ignored if tokenizer\_pool\_size is 0.

--tokenizer-pool-extra-config TOKENIZER\_POOL\_EXTRA\_CONFIG

Extra config for tokenizer pool. This should be a JSON string that will be parsed into a dictionary. Ignored if tokenizer\_pool\_size is 0.

--limit-mm-per-prompt LIMIT\_MM\_PER\_PROMPT

For each multimodal plugin, limit how many input instances to allow for each prompt. Expects a comma-separated list of items, e.g.: image=16,video=2 allows a maximum of 16 images and 2 videos per prompt. Defaults to 1 for each modality.

--enable-lora If True, enable handling of LoRA adapters.

--max-loras MAX\_LORAS

Max number of LoRAs in a single batch.

--max-lora-rank MAX\_LORA\_RANK

Max LoRA rank.

--lora-extra-vocab-size LORA\_EXTRA\_VOCAB\_SIZE

Maximum size of extra vocabulary that can be present in a LoRA adapter (added to the base model vocabulary).

--lora-dtype {auto,float16,bfloat16,float32}

Data type for LoRA. If auto, will default to base model dtype.

--long-lora-scaling-factors LONG\_LORA\_SCALING\_FACTORS

Specify multiple scaling factors (which can be different from base model scaling factor - see eg. Long LoRA) to allow for multiple LoRA adapters trained with those scaling factors to be used at the same time. If not specified, only adapters trained with the base model scaling factor are allowed.

--max-cpu-loras MAX\_CPU\_LORAS

Maximum number of LoRAs to store in CPU memory. Must be  $\geq$  than max\_num\_seqs. Defaults to max\_num\_seqs.

--fully-sharded-loras

By default, only half of the LoRA computation is sharded with tensor parallelism. Enabling this will use the fully sharded layers. At high sequence length, max rank or tensor parallel size, this is likely faster.

--enable-prompt-adapter

If True, enable handling of PromptAdapters.

--max-prompt-adapters MAX\_PROMPT\_ADAPTERS  
Max number of PromptAdapters in a batch.

--max-prompt-adapter-token MAX\_PROMPT\_ADAPTER\_TOKEN  
Max number of PromptAdapters tokens

--device {auto,cuda,neuron,cpu,openvino,tpu,xpu}  
Device type for vLLM execution.

--num-scheduler-steps NUM\_SCHEDULER\_STEPS  
Maximum number of forward steps per scheduler call.

--scheduler-delay-factor SCHEDULER\_DELAY\_FACTOR  
Apply a delay (of delay factor multiplied by previous prompt latency) before scheduling next prompt.

--enable-chunked-prefill [ENABLE\_CHUNKED\_PREFILL]  
If set, the prefill requests can be chunked based on the max\_num\_batched\_tokens.

--speculative-model SPECULATIVE\_MODEL  
The name of the draft model to be used in speculative decoding.

--speculative-model-quantization {aqlm,awq,deepspeedfp,tpu\_int8,fp8,fbgemm\_fp8,modelopt,marlin,gguf,gptq\_marlin\_24,gptq\_marlin,awq\_marlin,gptq,compressed-tensors,bitsandbytes,experts\_int8,qqq,neuron\_quant,None}  
Method used to quantize the weights of speculative model. If None, we first check the quantization\_config attribute in the model config file. If that is None, we assume the model weights are not quantized and use dtype to determine the data type of the weights.

--num-speculative-tokens NUM\_SPECULATIVE\_TOKENS  
The number of speculative tokens to sample from the draft model in speculative decoding.

--speculative-draft-tensor-parallel-size SPECULATIVE\_DRAFT\_TENSOR\_PARALLEL\_SIZE, -spec-draft-tp SPECULATIVE\_DRAFT\_TENSOR\_PARALLEL\_SIZE  
Number of tensor parallel replicas for the draft model in speculative decoding.

--speculative-max-model-len SPECULATIVE\_MAX\_MODEL\_LEN  
The maximum sequence length supported by the draft model. Sequences over this length will skip speculation.

--speculative-disable-by-batch-size SPECULATIVE\_DISABLE\_BY\_BATCH\_SIZE  
Disable speculative decoding for new incoming requests if the number of enqueue requests is larger than this value.

--ngram-prompt-lookup-max NGRAM\_PROMPT\_LOOKUP\_MAX  
Max size of window for ngram prompt lookup in speculative decoding.

--ngram-prompt-lookup-min NGRAM\_PROMPT\_LOOKUP\_MIN  
Min size of window for ngram prompt lookup in speculative decoding.

--spec-decoding-acceptance-method {rejection\_sampler,typical\_acceptance\_sampler}  
Specify the acceptance method to use during draft token verification in speculative decoding. Two types of acceptance routines are supported: 1) RejectionSampler which does not allow changing the acceptance rate of draft tokens, 2) TypicalAcceptanceSampler which is configurable, allowing for a higher acceptance rate at the cost of lower quality, and vice versa.

--typical-acceptance-sampler-posterior-threshold TYPICAL\_ACCEPTANCE\_SAMPLER\_POSTERIOR\_THRESHOLD  
Set the lower bound threshold for the posterior probability of a token to be accepted. T

his threshold is used by the TypicalAcceptanceSampler to make sampling decisions during speculative decoding. Defaults to 0.09

`--typical-acceptance-sampler-posterior-alpha` TYPICAL\_ACCEPTANCE\_SAMPLER\_POSTERIOR\_ALPHA  
A scaling factor for the entropy-based threshold for token acceptance in the TypicalAcceptanceSampler. Typically defaults to sqrt of `--typical-acceptance-sampler-posterior-threshold` i.e. 0.3

`--disable-logprobs-during-spec-decoding` [DISABLE\_LOGPROBS\_DURING\_SPEC\_DECODING]

If set to True, token log probabilities are not returned during speculative decoding. If set to False, log probabilities are returned according to the settings in SamplingParams. If not specified, it defaults to True.

Disabling log probabilities during speculative decoding reduces latency by skipping log prob calculation in proposal sampling, target sampling, and after accepted tokens are determined.

`--model-loader-extra-config` MODEL\_LOADER\_EXTRA\_CONFIG

Extra config for model loader. This will be passed to the model loader corresponding to the chosen load\_format. This should be a JSON string that will be parsed into a dictionary.

`--ignore-patterns` IGNORE\_PATTERNS

The pattern(s) to ignore when loading the model. Default to 'original/\*\*/\*' to avoid repeated loading of llama's checkpoints.

`--preemption-mode` PREEMPTION\_MODE

If 'recompute', the engine performs preemption by recomputing; If 'swap', the engine performs preemption by block swapping.

`--served-model-name` SERVED\_MODEL\_NAME [SERVED\_MODEL\_NAME ...]

The model name(s) used in the API. If multiple names are provided, the server will respond to any of the provided names. The model name in the model field of a response will be the first name in this list. If not

specified, the model name will be the same as the `--model` argument. Noted that this name(s) will also be used in model\_name tag content of prometheus metrics, if multiple names provided, metricstag will take the first one.

`--qlora-adapter-name-or-path` QLORA\_ADAPTER\_NAME\_OR\_PATH

Name or path of the QLoRA adapter.

`--otlp-traces-endpoint` OTLP\_TRACES\_ENDPOINT

Target URL to which OpenTelemetry traces will be sent.

`--collect-detailed-traces` COLLECT\_DETAILED\_TRACES

Valid choices are model,worker,all. It makes sense to set this only if `--otlp-traces-endpoint` is set. If set, it will collect detailed traces for the specified modules. This involves use of possibly costly and or blocking

operations and hence might have a performance impact.

`--disable-async-output-proc`

Disable async output processing. This may result in lower performance.

`--override-neuron-config` OVERRIDE\_NEURON\_CONFIG

override or set neuron device configuration.

`--lookahead-cache-config-dir` LOOKAHEAD\_CACHE\_CONFIG\_DIR

Folder path of lookahead cache config

`--cpu-decoding-memory-utilization` CPU\_DECODING\_MEMORY\_UTILIZATION

the memory is used for lookahead cache, which can range from 0 to 1. If unspecified, will use the default value of 0.15.

`--cpu-prefill-memory-utilization` CPU\_PREFILL\_MEMORY\_UTILIZATION

the memory is used for prefill cache, which can range from 0 to 1. If unspecified, will use the default value of 0.3.

--ignore-prompt-for-lookahead-cache

If True, the prompt will be ignored.

--enable-prefix-cache-offload

Enables prefix cache offloading

--apc-offload-not-lazy

If set, lazy launch of layer 2~n-1 will be disabled.

--apc-offload-min-access-threshold APC\_OFFLOAD\_MIN\_ACCESS\_THRESHOLD

Min threshold for evict offloading. Default 1.

--apc-offload-enable-hit-cnt

Enable hit count in APC.

--apc-offload-gpu-evictor-limit APC\_OFFLOAD\_GPU\_EVICTOR\_LIMIT

The free table size limited in gpu evictor. -1 default disable.

--disable-log-requests

Disable logging requests.

--max-log-len MAX\_LOG\_LEN

Max number of prompt characters or prompt ID numbers being printed in log. Default: Unlimited

除了兼容 vLLM 所有的配置参数外，TACO-LLM 还额外添加了以下参数配置：

# Lookahead-Cache

--lookahead-cache-config-dir LOOKAHEAD\_CACHE\_CONFIG\_DIR

Folder path of lookahead cache config

--ignore-prompt-for-lookahead-cache

If True, the prompt will be ignored.

--cpu-decoding-memory-utilization CPU\_DECODING\_MEMORY\_UTILIZATION

the memory is used for lookahead cache, which can range from 0 to 1. If unspecified, will use the default value of 0.15.

# Auto Prefix Cache CPU Offload

--enable-prefix-cache-offload

Enables prefix cache offloading

--cpu-prefill-memory-utilization CPU\_PREFILL\_MEMORY\_UTILIZATION

the memory is used for prefill cache, which can range from 0 to 1. If unspecified, will use the default value of 0.3.

--apc-offload-not-lazy

If set, lazy launch of layer 2~n-1 will be disabled.

# Sampling API

```
class taco_llm.SamplingParams(  
    n: int = 1,  
    best_of: Optional[int] = None,  
    presence_penalty: float = 0.0,  
    frequency_penalty: float = 0.0,  
    repetition_penalty: float = 1.0,  
    temperature: float = 1.0,  
    top_p: float = 1.0,  
    top_k: int = -1,  
    min_p: float = 0.0,  
    seed: Optional[int] = None,  
    use_beam_search: bool = False,  
    length_penalty: float = 1.0,  
    early_stopping: Union[bool, str] = False,  
    stop: Optional[Union[str, List[str]]] = None,  
    stop_token_ids: Optional[List[int]] = None,  
    ignore_eos: bool = False,  
    max_tokens: Optional[int] = 16,  
    min_tokens: int = 0,  
    logprobs: Optional[int] = None,  
    prompt_logprobs: Optional[int] = None,  
    detokenize: bool = True,  
    skip_special_tokens: bool = True,  
    spaces_between_special_tokens: bool = True,  
    logits_processors: Optional[Any] = None,  
    include_stop_str_in_output: bool = False,  
    truncate_prompt_tokens: Optional[Annotated[int, msgspec.Meta(ge=1)]] = None,  
    no_repeat_ngram_size: int = 0  
)  
    """Sampling parameters for text generation.
```

Overall, we follow the sampling parameters from the OpenAI text completion API (<https://platform.openai.com/docs/api-reference/completions/create>). In addition, we support beam search, which is not supported by OpenAI.

## Args:

- `n`: Number of output sequences to return for the given prompt.
- `best_of`: Number of output sequences that are generated from the prompt. From these `best_of` sequences, the top `n` sequences are returned. `best_of` must be greater than or equal to `n`. This is treated as the beam width when `use_beam_search` is True. By default, `best_of` is set to `n`.
- `presence_penalty`: Float that penalizes new tokens based on whether they

appear in the generated text so far. Values  $> 0$  encourage the model to use new tokens, while values  $< 0$  encourage the model to repeat tokens.

**frequency\_penalty:** Float that penalizes new tokens based on their frequency in the generated text so far. Values  $> 0$  encourage the model to use new tokens, while values  $< 0$  encourage the model to repeat tokens.

**repetition\_penalty:** Float that penalizes new tokens based on whether they appear in the prompt and the generated text so far. Values  $> 1$  encourage the model to use new tokens, while values  $< 1$  encourage the model to repeat tokens.

**temperature:** Float that controls the randomness of the sampling. Lower values make the model more deterministic, while higher values make the model more random. Zero means greedy sampling.

**top\_p:** Float that controls the cumulative probability of the top tokens to consider. Must be in  $(0, 1]$ . Set to 1 to consider all tokens.

**top\_k:** Integer that controls the number of top tokens to consider. Set to -1 to consider all tokens.

**min\_p:** Float that represents the minimum probability for a token to be considered, relative to the probability of the most likely token. Must be in  $[0, 1]$ . Set to 0 to disable this.

**seed:** Random seed to use for the generation.

**use\_beam\_search:** Whether to use beam search instead of sampling.

**length\_penalty:** Float that penalizes sequences based on their length. Used in beam search.

**early\_stopping:** Controls the stopping condition for beam search. It accepts the following values: ``True``, where the generation stops as soon as there are ``best_of`` complete candidates; ``False``, where an heuristic is applied and the generation stops when is it very unlikely to find better candidates; ``"never"``, where the beam search procedure only stops when there cannot be better candidates (canonical beam search algorithm).

**stop:** List of strings that stop the generation when they are generated. The returned output will not contain the stop strings.

**stop\_token\_ids:** List of tokens that stop the generation when they are generated. The returned output will contain the stop tokens unless the stop tokens are special tokens.

**include\_stop\_str\_in\_output:** Whether to include the stop strings in output text. Defaults to `False`.

**ignore\_eos:** Whether to ignore the EOS token and continue generating tokens after the EOS token is generated.

**max\_tokens:** Maximum number of tokens to generate per output sequence.

**min\_tokens:** Minimum number of tokens to generate per output sequence before EOS or `stop_token_ids` can be generated

**logprobs:** Number of log probabilities to return per output token.

When set to `None`, no probability is returned. If set to a non-`None` value, the result includes the log probabilities of the specified

number of most likely tokens, as well as the chosen tokens.

Note that the implementation follows the OpenAI API: The API will always return the log probability of the sampled token, so there may be up to `logprobs+1` elements in the response.

`prompt_logprobs`: Number of log probabilities to return per prompt token.

`detokenize`: Whether to detokenize the output. Defaults to True.

`skip_special_tokens`: Whether to skip special tokens in the output.

`spaces_between_special_tokens`: Whether to add spaces between special tokens in the output. Defaults to True.

`logits_processors`: List of functions that modify logits based on previously generated tokens, and optionally prompt tokens as a first argument.

`truncate_prompt_tokens`: If set to an integer `k`, will use only the last `k` tokens from the prompt (i.e., left truncation). Defaults to None (i.e., no truncation).

`no_repeat_ngram_size`:

If set to `int > 0`, all ngrams of that size can only occur once.

"""

除了兼容 vLLM 所有的采样参数外，TACO-LLM 还额外添加了以下采样参数：

`no_repeat_ngram_size`:

If set to `int > 0`, all ngrams of that size can only occur once.

# TACO LLM 性能

## 下载性能测试包

性能测试可以参见我们提供的性能测试包。

```
wget -c
```

## 解压

```
tar -zxvf taco_llm_demo.tar.gz  
cd taco_llm_demo
```

## 初始化环境

下载数据集、参考模型等。

```
#!/bin/bash
```

# 需要在有外网的环境下，先下载对应的测评需要用到的数据集 如果本目录有了，则忽略

```
wget -c https://taco-1251783334.cos.ap-shanghai.myqcloud.com/llm/data/ShareGPT_V3_unfiltered_cleaned_split.json
```

```
wget -c https://taco-1251783334.cos.ap-shanghai.myqcloud.com/llm/data/c4_sample.jsonl
```

```
wget -c https://taco-1251783334.cos.ap-shanghai.myqcloud.com/llm/data/medical_dialogue.json
```

```
wget -c https://taco-1251783334.cos.ap-shanghai.myqcloud.com/llm/data/github_sample.jsonl
```

# 下载参考模型Llama-2-7b-hf，注意更多模型可以在huggingface下载

```
wget -c https://taco-1251783334.cos.ap-shanghai.myqcloud.com/llm/llama/llama-2/Llama-2-7b-hf.tar
```

注意：

如果客户自己有数据集也可以在相关的脚本中调整、如果有相同的数据集，也可以不用再下载。

## 在线场景性能测试

运行 server 端

```
bash taco_bench_server.sh taco_llm
```

可以直接执行 server 端脚本，后面附加一个框架名称，例如 taco\_llm。启动该脚本的目的是创建一个server端的等待任务，待 client 请求处理。

server 脚本中关键参数：大多数参数可参照本文中的在线模型进行配置。以下是更多参数配置的说明：

```
chat_template="./llama.jinja"    # chat配置的模板路径，包里已经包含

# 设置prompt参数
SYSTEM_PROMPT_LENGTH=0
tgt_max_len=300                # 请求生成的最大长度
NUM_PROMPT_PRE_TGT=5          # 每个并发请求数
NUM_TGT=1                      # 每秒并发数

#设置服务器参数
host="127.0.0.1"              # 服务器地址
port="8007"                   # 服务端口
max_num_batched_tokens=10240  # 表示每次执行推理时支持最长的处理token数，多个batch一次处理的token总数。
max_num_seqs=32               # 后端支持最大的batch数
```

运行 client 端

```
bash taco_bench_client.sh taco_llm
```

client 脚本中关键参数：大多数参数可以参考本文中提到的在线模型。更多参数配置的说明如下：

```
DATASET_PATH="./ShareGPT_V3_unfiltered_cleaned_split.json"  # 数据集路径
MODEL_PATH="/models/Llama-2-7b-hf"                          # 模型路径
tp=1                                                         # 需要的GPU卡数量
TOKENIZER_PATH=$MODEL_PATH

# 设置prompt参数
SYSTEM_PROMPT_LENGTH=0
tgt_max_len=300                # 请求生成的最大长度
NUM_PROMPT_PRE_TGT=5          # 每个并发请求数
NUM_TGT=1                      # 每秒并发数

# 设置输出长度
output_len=100                # 设置每个请求最大输出数量

#设置服务器参数
backend=${1}                  # 服务端地址，需要和server配套
host="127.0.0.1"
port="8007"
```

```

ENABLE_PREFIX_CACHE=true      # true/false: 打开/关闭 Auto Prefix Cache功能
ENABLE_CACHE_OFFLOAD=true     # true/false: 打开/关闭 Cache Offload功能
ENABLE_HIT_CNT=true           # true/false: Cache命中情况打印
ENABLE_LOOKAHEAD=false        # true/false: 打开/关闭 lookahead功能
    
```

结果

我们将结果根据相关指标存储在本地的 results 目录中。

```

===== Serving Benchmark Result =====
Backend:                    taco_llm
Traffic request rate:      inf
Successful requests:       4
Benchmark duration (s):    9.58
Total input tokens:        1980
Total generated tokens:    1600
Total generated tokens (retokenized): 1602
Request throughput (req/s): 0.42
Input token throughput (tok/s): 206.76
Output token throughput (tok/s): 167.08
-----End-to-End Latency-----
Mean E2E Latency (ms):     2390.31
Median E2E Latency (ms):   2069.78
-----Time to First Token-----
Mean TTFT (ms):           656.38
Median TTFT (ms):         41.38
P99 TTFT (ms):            2427.64
-----Time per Output Token (excl. 1st token)-----
Mean TPOT (ms):           4.35
Median TPOT (ms):         4.51
P99 TPOT (ms):            5.50
-----Inter-token Latency-----
Mean ITL (ms):            12.43
Median ITL (ms):          12.38
P99 ITL (ms):             13.27
=====
    
```

taco\_llm\_demo/results/\*\*.csv

csv 表格结果如下：

Name	Backend	prefix	cache	offload	hit	lookahead	duration	input tokens	generated	throughput	throughput	E2E Latency	E2E Latency	TTFT (ms)	TTFT (ms)	TTFT (ms)	TPOT (ms)	TPOT (ms)	TPOT (ms)	ITL (ms)	ITL (ms)	ITL (ms)
20241220	taco_llm	false	false	false	true	9.576429	1980	1600	206.7577	167.0769	2390.311	2069.776	656.3841	41.38155	2427.637	4.345681	4.508368	5.503256	12.4275	12.37992	13.26822	

一键测试

客户也可以使用一个脚本完成 server/client 端测评部署并直接得到结果

```
bash taco_bench.sh taco_llm
```

注意：

实际上，这个脚本将上述的 server 端和 client 端两个脚本合并，在启动 server 端后等待10秒，然后启动 client 端。