

AIPerf: 大规模人工智能算力基准测试程序

关键词：人工智能 算力基准测试

陈文光 翟季冬 郑纬民 等
清华大学

背景

近年来，随着人工智能在自然语言处理、计算机视觉等领域上的快速发展以及在大规模算力上的普及，公众需要一个简单有效的指标来帮助判断系统的人工智能算力和整个高性能人工智能领域的发展状况。同时，一个好的指标也可以引领一个领域的健康持续发展。

然而，传统的高性能计算机评测方法和体系与当前人工智能需求的性能并不完全一致。例如，LINPACK 是一个目前被广泛采用的高性能计算机双精度浮点运算性能基准评测程序，国际超算 Top 500 榜单依据 LINPACK 值来进行排名，而典型的人工智能应用并不需要双精度浮点数运算。大部分人工智能训练任务以单精度浮点数或半精度浮点数为主，推理以 Int8 为主。

对大规模人工智能算力来说，制定一个简单有效的指标和测试方法并不是一件容易的事情。首先，大部分单个人工智能训练任务（例如训练一个推荐系统或者图像分类的神经网络模型）达不到全机上百张加速器卡规模的计算需求。很多人工智能应用，即使使用全机规模，其训练时间和准确率也可能没有改进。其次，如果要

测试规模变化的人工智能集群计算机，测试程序必须能够规模可变。我们必须明确，什么样的主流人工智能应用是规模可以任意调整的。最后，准确率的判定和计算是大规模人工智能算力评测与传统高性能计算基准评测之间的一个显著区别。是否需要使残差小于给定标准，是否要将准确度计入分数统计，同样是需要明确的问题。

目前，各大企业、高校和相关组织在人工智能性能基准测试领域已经有了很多探索，相继开发了各类基准评测程序，比如谷歌等公司主导的 MLPerf^[1]，小米公司的 MobileAI bench^[2]，百度公司的 DeepBench^[3]，中国人工智能产业发展联盟的 AIA DNN Benchmark^[4]，以及在双精度的 LINPACK 基础上改成混合精度的 HPL-AI^[5] 等。但是这些基准测试方案都不能很好地解决上述问题。根据 MLPerf 公

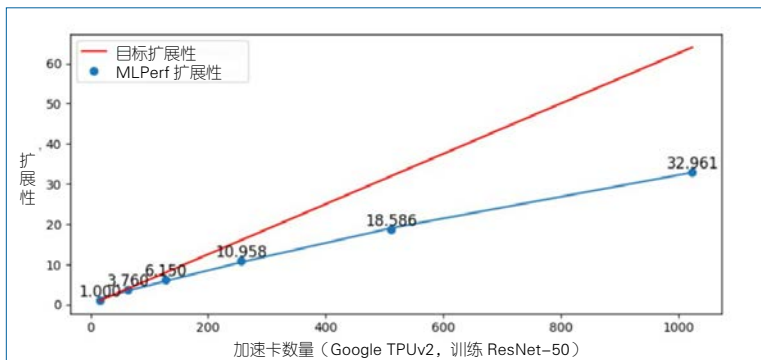


图1 MLPerf 扩展性瓶颈^[1]

公开发表的数据（如图1）可知，MLPerf程序在百张TPU加速卡以上规模测试下扩展性就会出现下滑，在千张TPU加速卡级别到达评测体系的扩展性瓶颈，该评测程序很难评价不同系统在该规模下人工智能算力的差异。

AIPerf 设计目标与思路

我们研制的人工智能算力基准测试程序 AIPerf，希望能满足如下目标：

1. 一个统一分数

基准测试程序应当报告一个分数作为被评测计算集群系统的评价指标。使用一个而不是多个分数能方便地对不同机器进行横向比较，以及方便对公众的宣传。除此之外，该分数应当随着人工智能计算集群的规模扩大而线性增长，从而能够准确评测不同系统规模下算力的差异。

2. 可变的问题规模

人工智能计算集群往往有着不同的系统规模，差异性体现在节点数量、加速器数量、加速器类型、内存大小等指标上。因此，为了适应各种规模的高性能计算集群，预期的人工智能基准测试程序应当能够通过变化问题的规模来适应集群规模的变化，从而充分利用人工智能计算集群的计算资源来体现其算力。

3. 具有实际的人工智能意义

具有人工智能意义的计算，例如神经网络运算，是人工智能基准测试程序与传统高性能计算机基准测试程序的重要区别，也是其能够检测集群人工智能算力的核心所在。人工智能基准测试程序应当基于当前流行的人工智能应用而构建。

4. 评测程序包含必要的多机通信

网络通信是人工智能计算集群设计的主要指标之一，也是其庞大计算能力的重要组成部分。面向高性能计算集群的人工智能基准测试程序应当包括必要的多机通信，从而将网络通信性能作为最终性能的影响因素之一。同时，基准测试程序中的多机通信模式应该具有典型的代表性。

针对以上设计目标，我们提出了基于自动化机

器学习（AutoML^[9]）来进行人工智能基准测试程序的设计。有如下三个原因：

1. AutoML 是通过算法自动搜索合适的神经网络模型结构，找到针对特定任务效果最好的解。因此该应用所需的计算资源极高，基础算法也包含训练模型本身，负载具有人工智能意义。

2. AutoML 目前包括超参数搜索和网络结构搜索两个方面。其中超参数搜索易于实现但搜索空间受到限制，而网络结构搜索则有着更大的搜索空间。通过结合这两种搜索方式，AutoML 能够充分利用大量的计算资源。

3. AutoML 具有足够的并行度，常常需要同时训练大量候选模型来对结构进行评估。同时，AutoML 搜索的结果虽然有一定的随机性，但整体上能找到的解的优劣程度随着搜索消耗的计算量增加而逐渐改善。

AIPerf 设计框架

AIPerf 基于微软 NNI 开源框架实现^[6]，以 AutoML 为负载，使用 Network Morphism 网络结构搜索和 TPE（Tree-structured Parzen Estimator）^[10] 超参搜索来寻找精度更高的神经网络结构和（或）超参数^[7]。

用户可以通过配置文件指定 AutoML 的相关参数，如训练使用的批大小（batch size）、最大 epoch 数、学习率、最大搜索模型总个数、最长搜索总时间、最大同时搜索模型个数（并发数）等多个参数。

通过配置文件指定 AutoML 任务后，用户可以在主节点上启动 AIPerf。如图2所示，启动后，主节点将持续地将 AutoML 任务以及当前最新的历史信息通过 SLURM 资源管理系统分发给各个工作节点。历史信息包括已经被工作节点完成的 AutoML 任务的网络模型及其精度。工作节点接收到任务与历史信息后，将根据历史信息开始搜索并生成一个新的网络模型。模型生成完成后，工作节点调用后端深度学习框架（默认使用 Keras 和 TensorFlow）进行神经网络训练。训练过程中，使用网络文件系统（Network File System, NFS）共享数据集文件，方便所有工作节点并发读取。

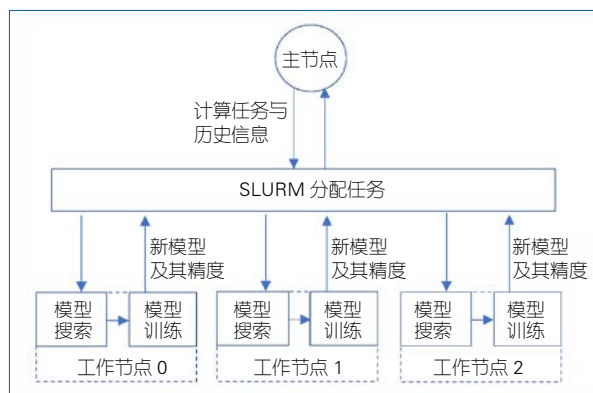


图2 AIPerf Benchmark 执行流程

工作节点根据一定策略（如训练 epoch 数，达到配置文件指定的最大 epoch 数或连续 10 个 epoch 没有精度提升）停止模型训练，并将此网络模型及其精度发送给主节点，主节点更新历史信息，并继续发送新的任务和历史信息。当任务完成数达到配置文件中指定的最大任务数量或 AIPerf 启动时间到达配置文件中指定的总时长时，AIPerf 结束。

在 AIPerf 执行过程中，所有工作节点上的任务都是异步执行的，即一个工作节点完成其任务时，不需要等待其他节点，可以直接将新网络模型及其精度发送回主节点，并从主节点接收下一个任务以及新的历史信息。

模型搜索与生成是 AutoML 中至关重要的一步，AIPerf 目前使用 Network Morphism 网络结构搜索和 TPE 超参搜索来寻找精度更高的神经网络结构和（或）超参数。Network Morphism 使用了网络形态技术来搜索和评估新的网络结构。在 AIPerf 中，工作节点首先会从父网络生成几个子网络，然后使用贝叶斯算法，从网络的历史训练结果来预测子网络的模型精度。接下来，选择预测模型精度最高的子网络进行训练。从父网络生成子网络包含一些层的变换，如变宽、变深，或在层中增加跳跃连接。TPE 是一种基于序列模型优化的方法，它能够根据历史模型精度来按顺序构造模型，估算超参的性能，随后基于此模型来选择新的超参数。

历史数据的数据量随着任务完成数增多而增多，且使用的工作节点越多，历史数据量的增长速

度也越快。由于模型搜索所需要的时间与历史数据量密切相关，当集群规模非常大的时候，AIPerf 执行后期可能需要花费非常多的时间用于模型搜索，这将导致得到的集群 AI 性能指标下降。

AIPerf 采用了多种方法来减少模型搜索所用时间。一种方法是当一个工作节点在进行模型训练的时候，其异构加速器被占用，而 CPU 基本是空闲的，此时 CPU 可以使用目前的历史数据提前进行模型搜索与生成，即使用模型训练的时间来掩盖模型搜索与生成的时间，其缺点是由于提前进行模型搜索，使用的历史数据并不是最新的，可能导致模型搜索效果较差；另一种方法是使用异构加速芯片，例如 GPU，加速模型搜索与生成的过程。模型搜索与生成需要对旧模型与新模型进行编辑距离的比较，从而找到预计最好的模型进行后续的实际训练。而传统编辑距离的计算由于依赖复杂难以并行，用异构加速芯片加速则是通过基于多边形变换的方法来消除复杂依赖，将编辑距离计算并行化，同时采用多边形拼接的方式来增大交给异构加速芯片的负载，从而充分利用加速芯片的计算资源。这个方法是直接对模型搜索生成过程进行加速，因此保证使用的历史数据是最新的。

AIPerf 评价指标

AIPerf 目前的评价指标是 Tops，即平均每秒处理的混合精度 AI 浮点操作数，也是本次 AIPerf 发榜的主要排名依据。评测指标的主要计算方法是，统计在规定的评测时间内所需的 AI 操作数，然后除以所需的评测时间，所得结果作为最后的评价指标。同时，搜索和训练的模型需要达到一定精度（默认为 70%），评价指标才被视为有效。

AIPerf 排行榜

2020 年 11 月 15 日，以“新算力，新基建，新经济”为主题的第二届中国超级算力大会（ChinaSC）在北京举行。大会上，基于 AIPerf 的国际人工智能能

性能算力排行榜首次发布^[8], 该榜单由中科院计算所研究员张云泉、清华大学教授陈文光、美国阿贡国家实验室研究员 Pavan Balaji 和瑞士苏黎世联邦理工大学教授

Torsten Hoefler

联合 ACM SIGHPC China 委员会共同发起。AIPerf 榜单前 10 机器如图 3 所示。

鹏城实验室研制的基于 ARM 架构和华为加速处理器的鹏城云脑二主机以 194527 Tops 的 AIPerf 算力荣登榜首, 其性能远超排名第二的联泰集群 NVIDIA 系统 (12 倍)。

鹏城云脑二主机包括 512 个计算节点, 4096 个华为 Ascend 910 (内存 32 GB) AI 加速卡, 采用 2048 个华为 Kunpeng 920 2.6 GHz CPU, 每个计算节点的物理内存为 2048 GB, 系统总内存为 1024 TB, 互连网络采用 100G RoCE, 操作系统为 EulerOS 2.8。

此外, 前 10 榜单还包括互联网视频公司、中软国际有限公司、清华同方内蒙古高性能服务中心、深圳早知道有限公司、中南大学、北京人工智能研究院的人工智能计算集群。

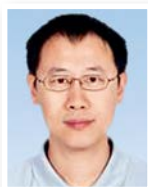
AIPerf 展望

AIPerf 基准评测程序还处于开发和完善阶段, 目前正在应用负载、硬件适应性以及国际推广等方面进行大力推进。AIPerf 目前只支持面向计算机视觉的人工智能应用程序。为了更好地评估大规模智能系统在各个典型应用领域的性能, 计划今后支持更多种网络搜索与训练算法和评测数据集, 同时支持更多类型的人工智能计算集群系统, 并积极扩大 AIPerf 在国际上

排名	研制厂商/单位	安装单位名称	AIPerf 值 (Tops)	系统名称/型号
1	鹏城实验室	鹏城实验室	194527.5	CloudBrainII, 512nodes, Ascend910(32GB) × 4096, Kunpeng 920 2.6GHz × 2048, RoCE 100G network
2	联泰集群	某视频公司	16361.28	2100IG Cluster, 276 nodes, NVIDIA Tesla T4(16GB) × 2208, Intel Xeon 5218 × 552, 10GbE/Intel network
3	联泰集群	某视频公司	14220.42	4800IG Cluster, 295 nodes, NVIDIA GTX 2080Ti × 2360, Intel Xeon E5-2680v4 × 590, 10GbE/Intel network
4	联泰集群	某科技公司	10360	4800IG Cluster, 70 nodes, NVIDIA Tesla V100(32G) × 560, Intel Xeon 6146 × 140, FDR IB/Mellanox network
5	中软国际科技服务有限公司	中软国际科技服务有限公司	5070	NVIDIA Tesla NVLink V100 Host, 32 nodes, NVIDIA Tesla NVLinkV100(32GB) × 256, Intel Xeon Skylake 6129 2.3GHz × 40, InfiniBand 100Gb/s network
6	同方股份有限公司	内蒙古高性能计算公共服务平台(青城之光)	3918.54	TANAC2020 Cluster, 72 nodes, NVIDIA Tesla PCIe V100(16GB) × 288, Intel Xeon Gold 6129 2.3GHz × 148, InfiniBand 100Gb/s network
7	深圳市早知道科技有限公司	深圳市早知道科技有限公司	1350	NVIDIA Tesla NVLink V100 Host, 8 nodes, NVIDIA Tesla NVLink V100(32GB) × 64, Intel Xeon Skylake 6151 3.0GHz × 40, InfiniBand 100Gb/s network
8	中南大学	中南大学	172.33	Central South University Platform, 1 node, NVIDIA Tesla V100(32GB) × 8, Intel Xeon Gold 6248 CPU 2.50GHz × 2, InfiniBand 100Gb/s network
9	某超级计算中心	某超级计算中心	147.65	NF, 5 nodes, NVIDIA Tesla V100(32GB) × 20, Intel Xeon Bronze 3106 1.70 GHz × 10, 10GbE/Intel network
10	英伟达	北京智源人工智能研究院	134.89	DGX 1, 1 node, NVIDIA Tesla V100(32GB) × 8, Intel Xeon CPU E5-2698 v4 2.20GHz × 2, Dual 10 GbE, 4 IB EDR network

图 3 2020 年 AIPerf 排行榜

的影响力, 希望能将其打造成具有国际影响力和公信力的大规模人工智能系统基准评测程序。



陈文光

CCF 副秘书长、理事、杰出演讲者, 曾任 CCF 编委、YOCSEF 主席 (2011~2012 年度)。清华大学教授, 兼任青海大学计算机系主任。主要研究方向为并行计算的编程模型、并行化编译和应用分析。cwg@tsinghua.edu.cn



翟季冬

CCF 高级会员, 2020 “CCF-IEEE CS 青年科学家奖”获得者。清华大学副教授。主要研究方向为高性能计算、并行程序性能分析和优化。zhaijidong@tsinghua.edu.cn



郑纬民

CCF 会士, CCF 前理事长 (2012—2016)。清华大学教授。中国工程院院士。主要研究方向为并行/分布处理、网络存储器等。zwm-dcs@tsinghua.edu.cn

其他作者: 任志祥 张云泉 余 腾 师天庵
谢 磊 钟闰鑫 刘永恒

参考文献

- [1] Mattson Peter, Reddi V J, Cheng C, et al. MLPerf: An industry standard benchmark suite for machine learning performance[J]. *IEEE Micro*, 2020, 40(2): 8-16.
- [2] Xiaomi. mobile-ai-bench[OL]. <https://github.com/XiaoMi/>

mobile-ai-bench.

- [3] Narang, S., and G. Diamos. "Baidu deepbench." (2017).
- [4] AIIA[OL]. <https://github.com/AIIABenchmark/AIIA-DNN-benchmark>
- [5] Dongarra J, Luszczek P, Tsai Y M. HPL-AI mixed-precision benchmark[OL].<https://icl.bitbucket.io/hpl-ai>.
- [6] Microsoft[OL]. <https://github.com/microsoft/nni>
- [7] Jin H, Song Q, Hu X. Auto-keras: An efficient neural architecture search system[C]// Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019.
- [8] AIPerf[OL]. www.aiperf.org
- [9] He X, Zhao K, Chu X. AutoML: A Survey of the state-of-the-art[J]. Knowledge-Based Systems, 2021, 212(5): 106622.
- [10] Bergstra J, Bardenet R, Bengio Y, et al. Algorithms for Hyper-Parameter Optimization[C]//Advances in Neural Information Processing Systems. 2011: 2546-2554.