云天励飞DESDK Benchmark工具使用手册

版本: v1.0.0

更新记录

版本	修改日期	修改说明		
V1.0.0	2020.12.9	初始版本		

目录

1概述

2 基本概念

- 2.1 算法精度和精度损失
- 2.2 硬件性能和端到端性能
- 2.3 时延
- 2.4 吞吐量

3 测试流程

- 3.1 准备模型和数据
- 3.2 源码编译
- 3.3 修改配置文件
- 3.4 开始测试

4 精度测试

- 4.1 分类模型精度测试
- 4.2 检测模型精度测试

5 性能测试

- 5.1 吞吐量测试
- 5.2 时延测试

1 概述

benchmark工具用于对模型在芯片上的精度和推理性能进行测试和统计,使用户对DeepEye芯片当前配置下的推理能力有定量的了解,帮助用户根据需要选择合理的配置参数。

本文档主要介绍如何使用benchmark工具对模型在DeepEye芯片上的推理的算法精度和性能进行测试和统计。

2基本概念

2.1 算法精度和精度损失

算法精度是指对模型推理的结果增加后处理后的结果与测试集已标注的结果对比得到的准确度。

对于分类模型来说后处理的结果一般是由分类概率从大到小排列的分类标签及置信度,验证时一般使用Top1和Top5,表示实际分类结果在测试集标注结果中命中第1个和前5个的概率。

对于检测模型来说一般是检测框的位置和置信度,验证时一般使用mAP(mean Average Precision),表示检测框与实际框的交并比IoU(Intersection over Union)的平均值。

将原始模型在框架下的运行结果和芯片上运行结果分别做算法精度,得到Top或mAP精度。两者相减可以得到芯片运行的精度损失,作为评价芯片运行精度的指标。

2.2 硬件性能和端到端性能

DeepEye1000芯片为协处理器,硬件性能是指在芯片侧直接运行模型的性能,端到端性能指芯片作为协处理时从主控侧运行模型的性能。从主控侧运行比从芯片侧直接运行要增加核间调度和传输时间。

两种性能一般都从时延和吞吐量两个维度衡量。benchmark结果中硬件时延和吞吐量被称为芯片时延和芯片吞吐,端到端时延和吞吐量被称为主控时延和主控吞吐。

当模型较大,单次推理时间较长,流水调度比较充分,调度和传输时间会被掩盖,吞吐量相差不大;当模型较小时,单次推理时间较短,流水调度不够充分,调度和传输时间掩盖不掉,吞吐量会有一定差别。

2.3 时延

时延是指单个推理任务从数据送入推理接口开始到获得推理结果之间所消耗的时间的平均值。一般来说,并行推理的任务越多,任务间冲突越频繁,单个任务的时延越长。因此,单个任务独立运行时时延为最小。

2.4 吞吐量

吞吐量是指单位时间内完成推理任务的次数。一般来说,并行推理会更加充分利用硬件资源。在一定的并行数量范围内,并行的任务越多,吞吐量越大。并行任务超过一定的数量,由于资源限制和调度开销,吞吐量会趋于平稳甚至下降。因此,可以根据并行任务数和吞吐量得到趋势图,进而推算出吞吐量最大值。

3 测试流程

3.1 准备模型和数据

精度测试需要准备标注后的测试集数据和标注数据,性能测试只需要准备测试集数据。

示例使用的测试集数据和标注数据放在/DEngine/data/datasets目录下。

3.2 源码编译

DEngine发布版本中包含了benchmark运行所需的主控和芯片侧库,用户可直接使用。同时提供了benchmark实现的源码供用户参考,如果需要重新编译源码,流程如下:

进入/DEngine/desdk/cpp目录,执行:

如果是 linux-x64 平台, 生成的主控侧库放在 /DEngine/desdk/platform/host_linux-x64/lib/libbenchmark.so。

参见《云天励飞DEngine集成开发指南》准备好csky芯片侧交叉编译环境,执行:

```
# sh build_device.sh
```

生成的芯片侧库放在/DEngine/desdk/platform/dev_linux-dp1000/lib/libbenchmark.so。

3.3 修改配置文件

benchmark执行脚本会自动解析同目录下benchmark.json,按照配置执行相应测试。

配置项意义:

- device_num 设备(芯片)数目,每个芯片执行相同测试,测试结果为总的吞吐统计
- test_suites 测试集数组,每个元素为一个测试集配置
- test_suite_name 测试集名称
- topn 分类模型精度测试,在推理结果前topn个结果中做匹配
- accuracy_datasets_path 精度测试数据集目录
- synset_path 模型分类结果文件路径
- val_path 精度测试标注文件路径
- fps_datasets_path 测试数据集路径
- models_root 待测试芯片模型根目录
- thread_list 线程数列表,如果配置为[1,2],则表示分别以1线程和2线程做两次测试,分别输出测试结果
- batch_list batchsize列表,遍历列表中的batchsize对模型做多次测试
- nnp_num_list 与batch_list相对应,每种batchsize下配置的nnp个数
- resource_num_list 与batch_list相对应,每种batchsize下配置的调度资源数
- num_criteria 测试执行推理次数,测试程序会遍历读取文件,不足num_criteria的部分会循环使用
- · accuracy_test 是否做精度测试
- fps_test 是否做吞吐量测试
- input_quesize 测试入口队列的深度,默认为32。此值如果太小会导致调度受限,硬件资源不能完全利用。
- test_models 待测试模型集合(字典)。

每一个键值对以模型子目录名作为key, do_test表示此模型是否执行测试; shape表示此模型的输入 [h,w]; rgb表示模型输入格式; nnp_num_list和resource_num_list也可以在这里覆盖配置, 仅对本模型 生效。

配置文件示例内容如下:

```
{
   "device_num":1,
   "test_suites":
```

```
[{
    "test_suite_name":"normal",
    "accuracy_datasets_path":"/DEngine/data/datasets/ILSVRC2012",
    "synset_path":"/DEngine/data/datasets/synset_1000.txt",
    "val_path":"/DEngine/data/datasets/ILSVRC2012_val.txt",
    "fps_datasets_path":"/DEngine/data/datasets/ILSVRC2012",
    "models_root":"/DEngine/model/dp1000/caffe_squeezenet_v1.1",
    "thread_list":[1],
    "batch_list":[1,2,4],
    "nnp_num_list":[1,2,4],
    "resource_num_list":[4,4,4],
    "topn":5,
    "time_criteria": 10.0,
    "num_criteria": 50,
    "accuracy_test": false,
    "fps_test": true,
    "input_quesize": 32,
    "test_models":
      "caffe_squeezenet_v1.1": {"do_test":true, "shape":[227,227], "rgb":true,
"resource_num_list":[8,8,8]}
 }]
}
```

3.4 开始测试

进入/DEngine/desdk/tools/benchmark目录,执行命令:

```
# sh /DEngine/run.sh benchmark.py
```

测试完成后,将生成benchmark.xlsx文件,内容如下图所示。

機型の型	模型名称	MACK	特度		ISER	Batch	m	RES	manager c	THE WALLS	iantiants (**)	(62568/Dat (625)	
			tepl	tep5	nAP	35/66/80	Batch	***	Ans	銀件利用 (65)	新社的紅((20)	MATERIAL MAY	MARKAGE (1957)
							1	1	*	314	65	1192.032	65,36001633
							2	2	8	324	331	T45.384	130.336284
	ceffe_rematio	115/8/2012				1	4	4	8	418	211	548, 418	207. 9170193
							1	1	4	754	24	3184.348	24.71222013
							2	2	4	830	46	2152, 23	45.3535966
	outle_resuet50	TLSVBC2012				1	4	4	4	1297	63	2219.322	62.49103244
							1	1	4	1794	10	1990.64	10.4970461T
							2	2	4	1941	20	9089.624	19.04678831
	caffa_resnet152	ILSM8C2012				1	4	4	4	3027	27	540T, S2	26, 7645284
							1	1		126	157	453.16	196, 2614153
							z	z		133	362	192, 495	307, 3327023
	caffs_mebilemetryl	ELSVECZ012				1	4	4		357	521	50.43	345, 1457044
							1	1		127	157	450, 450	155, 1459451
							2	2		131	364	234.90	299, 2355991
	eafife_nobiliones=v0	TLSVBC2012				i	4	4		295	446	65.53	347. 644T34T
							1	1	4	1743	10	T149.96	11.0002002
							2	2	4	1838	21	46T3, T88	21, 323321
	caffs_inception_v3	TLENGC2012				1	4	4.	4	2232	36	3844, 228	36,58643388
							1	1	4	2130	8	8792,148	9.004723462
							2	2	4	2235	17	5870.838	17.16255399
	cuffs_inneption_v4	DISMICROID				1	4	4	4	2722	30	4829.902	29.70392768
							1	1		200	90	174.892	98.0006/501
							2	2		214	196	471. 122	194, 5202799
	eaffa_squeeemat_vt.0	TLEWRC2012				i	4	4.	8	290	364	414.62	294, 359196T
							1	1	*	102	196	348, 946	193, 3448LTB
							2	2		130	315	174.382	305, 3051708
	caffs agreement vi. 1	ETEMBES015				1	4	4.	8	294	449	50.712	340, 2558250
							1	1	4	2852	6	11830.996	6. T353T5248
							- 1	2	4	3015	13	7836,088	12.9400836T
	oaffe_rggif	TLSVBC2012				1	4	4	4	2002	21	6715.282	21.59134729
	2						1	1	4	439	44	1700.76	43.97736796
							2	2	4	127	54	1993, 594	52.9495449T
	caffa_proplanet	TLSV8C2012				1	4	4	4	1343	60	2318.926	61, 31039968

4 精度测试

精度测试是指使用指定数据集和标注结果对指定模型进行算法精度测试,以评估芯片运行的精度损失。

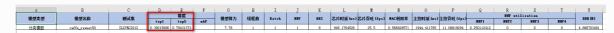
精度测试需要配置的key有accuracy_test(true)、accuracy_datasets_path、val_path。测试时使用的线程、batch等参数不影响结果,都用1即可。

精度测试分为分类模型精度测试和检测模型精度测试。分类模型精度测试目标为测试TOP1、TOP5的分类精度;检测模型的精度测试目标为mAP精度。

4.1 分类模型精度测试

分类模型精度测试时,需要准备好数据集文件 synset_1000.txt 、标签文件 ILSVRC2012_val.txt ,然后在benchmark.json文件中的synset_path和val_path字段中分别指定这两个文件,并设置accuracy_test为true。

精度测试结果在benchmark.xlsx文件中体现,如下图所示。



4.2 检测模型精度测试

注:目前检测模型精度测试功能正在完善中,暂不支持。

5 性能测试

性能测试是指使用指定数据集对指定模型进行多种模式和参数下的性能测试,以评估芯片运行的效率。 性能测试主要分为吞吐量测试和时延测试两种。

5.1 吞吐量测试

吞吐量测试模式下,测试程序将所有测试数据一次性发给芯片,驱动芯片以该配置下的最大能力执行, 得到最大的吞吐量数据。

由于芯片同一时间处理能力有限,未及时处理的测试数据会暂时存放在测试入口队列,等待前一批数据处理完之后进行处理。而在测试入口队列中的等待时间也会被记入时延,导致吞吐量测试模式下时延数据虚高,因此吞吐量测试模式下不用关注时延数据。

测试入口队列深度(input_quesize)减小会降低时延,但太小会导致硬件资源不能充分调度,导致吞吐量下降。

通过配置线程数、batch、nnp、res的不同组合来进行测试,可以获取一个最优配置达到最优的吞吐性能。一般情况,小模型可以适当配置res高一些提高性能,大模型配置多线程以提高吞吐。未达到最高性能应将nnp尽量配满,但多nnp情况下芯片mac利用率可能会有下降。

5.2 时延测试

时延测试模式下,测试程序将测试数据按照一定规则发给芯片,得到该规则下数据的最大、最小和平均 时延。

数据的发送规则一般为单一发送、平均间隔发送和泊松分布发送。

注:目前时延测试功能正在完善中,暂不支持。