



# GPU云服务器 最佳实践

文档版本: 20220602



# 法律声明

阿里云提醒您在阅读或使用本文档之前仔细阅读、充分理解本法律声明各条款的内容。 如果您阅读或使用本文档,您的阅读或使用行为将被视为对本声明全部内容的认可。

- 您应当通过阿里云网站或阿里云提供的其他授权通道下载、获取本文档,且仅能用 于自身的合法合规的业务活动。本文档的内容视为阿里云的保密信息,您应当严格 遵守保密义务;未经阿里云事先书面同意,您不得向任何第三方披露本手册内容或 提供给任何第三方使用。
- 未经阿里云事先书面许可,任何单位、公司或个人不得擅自摘抄、翻译、复制本文 档内容的部分或全部,不得以任何方式或途径进行传播和宣传。
- 由于产品版本升级、调整或其他原因,本文档内容有可能变更。阿里云保留在没有 任何通知或者提示下对本文档的内容进行修改的权利,并在阿里云授权通道中不时 发布更新后的用户文档。您应当实时关注用户文档的版本变更并通过阿里云授权渠 道下载、获取最新版的用户文档。
- 4. 本文档仅作为用户使用阿里云产品及服务的参考性指引,阿里云以产品及服务的"现状"、"有缺陷"和"当前功能"的状态提供本文档。阿里云在现有技术的基础上尽最大努力提供相应的介绍及操作指引,但阿里云在此明确声明对本文档内容的准确性、完整性、适用性、可靠性等不作任何明示或暗示的保证。任何单位、公司或个人因为下载、使用或信赖本文档而发生任何差错或经济损失的,阿里云不承担任何法律责任。在任何情况下,阿里云均不对任何间接性、后果性、惩戒性、偶然性、特殊性或刑罚性的损害,包括用户使用或信赖本文档而遭受的利润损失,承担责任(即使阿里云已被告知该等损失的可能性)。
- 5. 阿里云网站上所有内容,包括但不限于著作、产品、图片、档案、资讯、资料、网站架构、网站画面的安排、网页设计,均由阿里云和/或其关联公司依法拥有其知识产权,包括但不限于商标权、专利权、著作权、商业秘密等。非经阿里云和/或其关联公司书面同意,任何人不得擅自使用、修改、复制、公开传播、改变、散布、发行或公开发表阿里云网站、产品程序或内容。此外,未经阿里云事先书面同意,任何人不得为了任何营销、广告、促销或其他目的使用、公布或复制阿里云的名称(包括但不限于单独为或以组合形式包含"阿里云"、"Aliyun"、"万网"等阿里云和/或其关联公司品牌,上述品牌的附属标志及图案或任何类似公司名称、商号、商标、产品或服务名称、域名、图案标示、标志、标识或通过特定描述使第三方能够识别阿里云和/或其关联公司)。
- 6. 如若发现本文档存在任何错误,请与阿里云取得直接联系。

# 通用约定

格式	说明	样例
⚠ 危险	该类警示信息将导致系统重大变更甚至故 障,或者导致人身伤害等结果。	⚠ 危险 重置操作将丢失用户配置数据。
⚠ 警告	该类警示信息可能会导致系统重大变更甚 至故障,或者导致人身伤害等结果。	會告 重启操作将导致业务中断,恢复业务 时间约十分钟。
〔∫〉 注意	用于警示信息、补充说明等,是用户必须 了解的内容。	大意 权重设置为0,该服务器不会再接受新 请求。
? 说明	用于补充说明、最佳实践、窍门等 <i>,</i> 不是 用户必须了解的内容。	⑦ 说明 您也可以通过按Ctrl+A选中全部文件。
>	多级菜单递进。	单击设置> 网络> 设置网络类型。
粗体	表示按键、菜单、页面名称等UI元素。	在 <b>结果确认</b> 页面,单击 <b>确定</b> 。
Courier字体	命令或代码。	执行    cd /d C:/window    命令 <i>,</i> 进入 Windows系统文件夹。
斜体	表示参数、变量。	bae log listinstanceid
[] 或者 [alb]	表示可选项,至多选择一个。	ipconfig [-all -t]
{} 或者 {alb}	表示必选项,至多选择一个。	switch {act ive st and}

# 目录

1.搭建GRID驱动的License Server	05
1.1. 搭建Linux环境License Server	05
1.2. 搭建Windows环境的License Server	11
2.在GPU实例上部署NGC环境	21
3.GPU AI模型训练最佳实践	25
4.在GPU实例上使用RAPIDS加速机器学习任务	27
5.在GPU实例上使用RAPIDS加速图像搜索任务	37
6.RAPIDS加速机器学习最佳实践	45
7.RAPIDS加速图像搜索最佳实践	47

# 1.搭建GRID驱动的License Server

# 1.1. 搭建Linux环境License Server

使用GRID驱动必须申请相应的License,同时需要将License部署在License Server上,本文为您介绍搭建Linux 操作系统的License Server的具体操作。

# 背景信息

本文以Ubuntu 18.04操作系统为例,如果您使用了其他Linux操作系统,如CentOS,需要根据实际情况替换部分命令。更多信息,请参见NVIDIA官方文档。

NVIDIA GRID License的工作原理图如下所示:



您需要从NVIDIA License中心获取License文件,并将License文件部署在已搭建的License Server上。然后使用 您的GPU或vGPU实例,通过网络访问License Server激活GRID License。

## 步骤一: 创建ECS实例

创建一台Linux操作系统的ECS实例作为License Server。具体创建操作,请参见使用向导创建实例。

配置项	说明
实例规格	请选择配置高于2 vCPU 4 GB的ECS实例规格。如果您的业务需要大量的License,建议您 至少选择配置高于4 vCPU 16 GB的ECS实例规格,从而获得更高的性能及稳定性。
镜像	选择Linux操作系统镜像。本文以Ubuntu 18.04为例。
存储	请选择40 GiB以上系统盘。
带宽	建议您选择按使用流量计费,并将带宽峰值设置为100 Mbps。

您需要关注如下配置项,其他配置项根据自身业务场景设置即可。

## 步骤二:安装Java运行环境

1. 远程连接您用于搭建License Server的ECS实例。

具体连接操作,请参见<mark>通过密码或密钥认证登录Linux实例</mark>。

2. 运行以下命令,验证当前实例是否已安装Java运行环境。

java -version

#### 如果回显类似如下Java版本信息,表示已安装Java运行环境。否则,请继续执行后续操作完成安装。

java version "1.7.0\_51" OpenJDK Runtime Environment (rhel-2.4.5.5.el7-x86\_64 u51-b31) OpenJDK 64-Bit Server VM (build 24.51-b03, mixed mode)

3. 安装OpenJDK。

sudo apt install default-jdk

⑦ 说明 如果提示 Unable to locate package default-jdk 信息,请先执行 apt update 命 令后,再执行以上命令进行安装。

回显类似如下信息,表示OpenJDK已安装成功。

```
...
Running hooks in /etc/ca-certificates/update.d...
done.
done.
```

## 步骤三:安装并运行Apache Tomcat服务器

1. 使用Linux发行版的软件包管理器安装所需的Apache Tomcat软件包。

sudo apt install tomcat8

2. 安装完成后,运行以下命令,将Tomcat服务设置为开机自启动。

sudo systemctl enable tomcat8.service

3. 启动Tomcat服务。

sudo systemctl start tomcat8.service

- 4. 通过本地Web浏览器访问http://localhost:8080, 验证Tomcat服务是否可用。
  - ⑦ 说明 请将访问地址中的*localhost*替换为您的ECS实例的公网地址。

#### 如果显示如下内容,表示Tomcat服务已安装完成。

#### It works !

```
If you're seeing this page via a web browser, it means you've setup Tomcat successfully. Congratulations!
This is the default Tomcat home page. It can be found on the local filesystem at /ver/lib/towast8/wakeps/2007/index.html.
Tomcat8 Veterans might be pleased to learn that this system instance of Tomcat is installed with catultage in /ver/lib/towast8 and catultage in /ver/lib/towast8 and catultage in /ver/lib/towast8 and catultage in /ver/lib/towast8.
You might consider installing the following packages, if you haven't already done so:
tomcat8-adocs: This package installs a web application that allows to browse the Tomcat 8 documentation locally. Once installed, you can access it by clicking here.
tomcat8-adom: This package installs a web application that allows to access the Tomcat 8 documentation locally. Once installed, you can access it by clicking here.
tomcat8-adom: This package installs a web application that allows to access the Tomcat 8 documentation locally. Once installed, you can access it by clicking here.
tomcat8-adom: This package installs two web applications that can help managing this Tomcat instance. Once installed, you can access the manager webapp and the host-manager webapp.
NOTE: For security reasons, using the manager webapp is restricted to users with role "manager.ui". The host-manager webapp is restricted to users with role "manager.ui".
```

# 步骤四:安装License Server软件

您需要准备提前获取到的setup.bin安装文件或该文件所在的.zip压缩包。

- 1. 下载setup.bin安装程序或解压该程序所在的.zip压缩包。
  - 。 如果您通过NVIDIA官网获取了License Server的安装包,请运行以下命令,解压安装包。

unzip **安装程序.**zip

? 说明 请替换*安装程序.zip*为您获取到的安装包名称。

○ 如果您通过阿里云获取了License Server安装程序的下载地址,请运行以下命令,下载安装程序。

wget https://grid-9-4-zyy.oss-cn-hangzhou.aliyuncs.com/setup.bin

2. 运行以下命令,为安装程序添加执行权限。

chmod +x setup.bin

- 3. 安装License Server软件。
  - i. 以root用户运行安装程序。

sudo ./setup.bin -i console

ii. 在Introduction部分, 单击回车继续。



iii. 在License Agreement部分,请通过每一次单击回车进行翻页并接受许可协议。

当您达成许可协议时,系统会提示您接受许可协议条款,请输入Y,并单击回车。

DO YOU ACCEPT THE TERMS OF THIS LICENSE AGREEMENT? (Y/N): Y

- iv. 在Choose Install Folder部分,请单击回车,保持默认的License Server软件安装路径。
- v. 在**Choose Local Tomcat Server Path**部分, 输入Tomcat的本地路径, 默认为/*var/lib/tomcat 版本号*, 例如: /*var/lib/tomcat8*。
- vi. 在Choose Firewall Options部分,确认需要在防火墙中打开的端口,单击回车,保持默认选项即可。



vii. 在Pre-Installation Summary部分,确认信息并单击回车启动安装。



viii. 在Install Complete部分,单击回车,结束安装。

Install Complete
License Server has been successfully installed to:
/opt/flexnetls/nvidia
PRESS <enter> TO EXIT THE INSTALLER:</enter>

## 步骤五:在NVIDIA License上创建License Server

1. 前往NVIDIA License,并使用您申请License的邮箱登录。

冬 NVIDIA. LICENSING	Login	
€ ENTERPRISE SUPPORT	E-mail	
		LOG IN

2. 在Dashboard页面,单击License Servers区域下CREATE LICENSE SERVER按钮。

€ → C a uu	icensing.rvidia.com					ê 🕸 🕸 🔕
	ENSING   Dashboard				MIDIA Application Hub	and the second s
🔂 DASHBOARD	Entitlements		MANAGE ENTITLEMENTS	License Servers		MANAGE SERVERS CREATE SERVER
ENTITLEMENTS	» ENTITLEMENT (PAK ID) / FEATURE □	START DATE EXPRATION	ALLOCATED / TOTAL	» LICENSE SERVER	BOUND SERVICE INSTANCE	SERVICE INSTANCE TYPE
LUCENSE SERVE     SOFTWARE DO     VIRTUAL GROU     VIRTUAL GROU     SERVICE INSTAN     LEASES     EVENTS     AS USER MANAGEB     USER MANAGEB	s 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2			You	do not have any locate strong, would you	like to create one?
C ENTERPRISE SU	PPORT					

Entitlements区域下,是您目前拥有的全部License。

3. 在打开的Create License Server页面,填写相关信息,然后单击CREATE LICENSE SERVER。

#### GPU云服务器

Server Name	Feature	Lice	nses
Name this license server	Select a feature	✓ 1	
Description	Added Features		
Provide a short description	FEATURE		
	No featu	res have been added ye	t
MAC Address			
MAC Address (XX:XX:XX:XX:XX:XX or XX-XX-XX-XX-XX-XX)			
Failover server configuration is optional. If configuring, you must provide a name AND MAC address			
Failover License Server			
Failover License Server			
Failover MAC Address			
Falley as MAC Address			

#### 必填参数说明如下:

参数	说明
Server Name	自定义您需要的License Server名称。
MAC Address	填写您用于搭建License Server的ECS实例的MAC地址。 您可以登录实例,使用 ipconfig -a 命令进行查询。
Feature	根据需要选择,并输入待添加的License数量,单击 <b>ADD</b> 。

4. 创建完成后,在License Servers页面单击 \_\_\_\_按钮,下载License文件。

Li	cense Serve	rs	MANAGE SERVERS CREATE SERVER
≈	LICENSE SERVER	BOUND SERVICE INSTANCE	SERVICE INSTANCE TYPE
~	test LEGACY	n/a	n/a
		⊻ ∕ ⊝	⊕ Ø <b>■</b>
	FEATURE		TOTAL
	NVIDIA RTX Virtual	Workstation-5.0	10

# 步骤六:导入License文件

1. 通过本地Web浏览器访问http://localhost:8080/licserver, 进入License Server管理界面。

? i	<b>シ</b> 明 ù	, 与将访问地址中的 <i>localhost</i> 替换为您的ECS实例的公网地址。
-----	--------------	---

ועם 🧇	DIA.		
	Licensed Clients		
License Server	Licensed Clients with features consumed or reser	ved. Click a Client ID for further details.	
<ul> <li>Essenations</li> <li>Linemat Fasture Linema</li> </ul>	Olert ID	Client ID Type	Client Type
License Menagement	Total number of records: 0		
<ul> <li>Legin</li> </ul>			Copyright (c) 2003 MVIDIA Corporation: A3 Rights Reserved: 2003 05 0 2816036
Ucerse Client Narager	4		
<ul> <li>Sections</li> </ul>			

- 2. 在左侧导航栏的License Server区域,单击License Management。
- 3. 在License Management页面,单击Upload license file (.bin file)后的选择文件,选择本地的

### License文件,然后单击打开(O)。

× • • • • • • • • • • • • • • • • • • •
↓         ↓

4. 单击Upload。

页面提示如下信息,表示License文件已导入成功。

License Management	
Successfully applied license file to license server.	
Briwse fur the License file you received from the NMDA Licensing portal, and then click Upload to process the License file.	
• Upload license file ( bin file) 建造存任何文件	
	Cancel Upload

您可以在左侧导航栏的License Server区域,单击License Feature Usage,查看License数量以及使用情况。

Sareh (san-arrathal Sareh by Fasture V)	
Heatures	1 😰
Olick the feature table header row to obtain a single sorted non-paginated list. Features were fit or researced for clients. Click, electron name for usage details.	
- Futurn Winson Count Available Expliy Windor String	
<u>later and 10 a Dure and 10 D</u>	_

# 步骤七:测试网络连接和访问

以下操作以创建一台Windows操作系统的GPU虚拟化实例vgn6i为例,您也可以直接应用您已创建的GPU实例。

1. 创建一台GPU实例。

具体操作,请参见创建未配备驱动的GPU实例。

本示例选择的实例规格为GPU虚拟化实例vgn6i,操作系统为Windows Server 2019 数据中心版 64位中 文版。

2. 安装GRID驱动,然后添加License服务器并激活License。具体操作,请参见在GPU实例上安装GRID驱动 (Windows)。

激活License时,请填写您搭建了License Server的ECS实例公网IP,端口号为7070。

🛃 NVIDIA 控制面板	
文件(F) 编辑(E) 桌面(K) 帮助(H)	
🕲 fil 🔹 🕑 🐇	
选择—顶任务 □-3D 设置 □-通过预览调整图像设置	管理许可证
□ □ · 忙管理 30 设置 □ - 许可 □ □ □ <mark>□ □ 管理许可证</mark>	您可以通过应用许可证来启用附加功能。
	许可证类型: 夕您的系统已获 Quadro Virtual Data Center Workstation 许可。
	一级许可证服务器:
	9两山方: [7070
	二级许可证服务器:
	滅口号:
	说明: 主要许可证服务器聆听许可请求的演口编号。默认是 7070。

3. 前往License Server管理界面(*http://localhost:8080/licserver*),在左侧导航栏的License Server区域,单击Licensed Clients,即可查看到GPU实例信息,表示该NVIDIA vGPU软件客户端虚拟机目前正在使用License的功能。

Licensed Clients				
Licensed Clients with features consumed or reserved. Cl	ick a Client ID for further details.			
Client ID Page 1 of 1 Go to page 1 v Total number of records: 1	E	Client ID Type		Olient Type VIRTUAL
单击客户端ID,可	「以查看该客户端的)	羊细信息	l.	
Client Details				
Client Details				
Client ID: 0 Device ID Type: ETHERNET Client Type: VIRTUAL Client Expiry: 2022				
Licensed Features				
Feature Name Quadro	Version 5.0	Used 1	Expiry 20	Vendor String Quad

# 1.2. 搭建Windows环境的License Server

使用GRID驱动必须申请相应的License,同时需要将License部署在License Server上,本文为您介绍搭建Windows操作系统的License Server的具体操作。

# 背景信息

NVIDIA GRID License的工作原理图如下所示:



您需要从NVIDIA License中心获取License文件,并将License文件部署在已搭建的License Server上。然后使用 您的GPU或vGPU实例,通过网络访问License Server激活GRID License。

# 步骤一: 创建ECS实例

创建一台Windows操作系统的ECS实例作为License Server。具体创建操作,请参见使用向导创建实例。

您需要关注如下配置项,其他配置项根据自身业务场景设置即可。

配置项	说明
实例规格	请选择配置高于2 vCPU 4 GB的ECS实例规格。如果您的业务需要大量的License,建议您 至少选择配置高于4 vCPU 16 GB的ECS实例规格,从而获得更高的性能及稳定性。
镜像	选择Windows操作系统镜像。本文以Windows Server 2019 数据中心版 64位中文版为 例。
存储	请选择40 GiB以上系统盘。
带宽	建议您选择按使用流量计费,并将带宽峰值设置为100 Mbps。

## 步骤二:安装Java运行环境

1. 远程连接您用于搭建License Server的ECS实例。

具体连接操作,请参见通过密码或密钥认证登录Windows实例。

- 2. 请前往ojdkbuild下载OpenJDK JRE安装包。
- 3. 安装JRE。

🖟 OpenJDK 1.8.0_322-1-c	jdkbuild Setup	2_3		×
	Completed the Op 1.8.0_322-1-ojdkl	enJDK build Setup Wiz	ard	
	Click the Finish button to e:	vit the Setup Wizard.		
	Back	Finish	Cance	2]

4. 新建系统变量JAVA\_HOME,并将取值设置为JRE的jre文件夹所在的绝对路径。

例如,将路径设置为C:\Program Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.201-1\jre。

Path	C:\Users\Administrator\AppData\Local\Microsoft\WindowsA
TEMP	C:\Users\Administrator\AppData\Local\Temp
ТМР	C:\Users\Administrator\AppData\Local\Temp
	新建(N) 编辑(E) 删除(D)
统变量(S)	
变量	值
ComSpec	C:\Windows\system32\cmd.exe
DriverData	C:\Windows\System32\Drivers\DriverData
JAVA_HOME	C:\Program Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre
NOWBER_OF_PROCE	SSORS 4
US Dath	Windows_NI
PATHEXT	.COM; EXE; BAT; CMD; VBS; VBE; JS; JSE; WSF; WSH; MSC

- ⑦ 说明 请确保路径不包含任何尾随字符,例如反斜线(\)或空格。
- 在系统变量Path中,确认是否已存在java.exe程序所在的绝对路径。
   通常情况下,安装IRE时,会自动添加此路径。

emRoot%\system32 emRoot% emRoot% EMROOT%\System32\WindowsPowerShell\v1.0\ EMROOT%\System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	新建(N) 编辑(E) 浏览(B) 删除(D)
emRoot%\system32 emRoot% emRoot%\System32\Wbem 'EMROOT%\System32\WindowsPowerShell\v1.0\ 'EMROOT%\System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	新建(N) 编辑(E) 浏览(B) 删除(D)
emRoot%(system32 emRoot%(System32\Wbem 'EMROOT%(System32\WindowsPowerShell\v1.0\ 'EMROOT%(System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jire\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jire\bin	新建(N) 编辑(E) 浏览(B) 删除(D)
emRoot% emRoot%\System32\Wbem 'EMROOT%\System32\WindowsPowerShell\v1.0\ 'EMROOT%\System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	编辑(E) 浏览(B) 删除(D)
emRoot%(System32(Wbem 'EMROOT%(System32(WindowsPowerShell\v1.0\ 'EMROOT%(System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	编辑(E) 浏览(B) 删除(D)
TEMROOT%(System32\WindowsPowerShell\v1.0\ TEMROOT%\System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	浏览(B) 删除(D)
EMROOT%\System32\OpenSSH\ gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	浏览(B) 删除(D)
gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\bin gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	删除(D)
gram Files\ojdkbuild\java-1.8.0-openjdk-1.8.0.322-1\jre\bin	删除(D)
	上我们的
	下移(0)
	编辑文本(T)
	确定

# 步骤三:安装License Server软件

您需要准备提前获取到的License Server安装程序。

- 1. 解压License Server的.zip文件,并运行setup.exe安装程序。
- 2. 在Introduction页面, 单击Next。



3. 在License Agreement页面,选中I accept the terms of the License Agreement,然后单击Next。



4. 在Apache License Agreement页面,选中I accept the terms of the License Agreement,然 后单击Next。



5. 在Choose Install Folder页面,设置目标地址,然后单击Next。



6. 在Choose Firewall Options页面,保持默认选项License server (port 7070),然后单击Next。



7. 在Pre-Installation Summary页面,确认配置信息,单击Install。



8. 安装完成后,单击Done。



步骤四:在NVIDIA License上创建License Server

1. 前往NVIDIA License,并使用您申请License的邮箱登录。

📀 NVIDIA. LICENSING   l	_ogin
€ ENTERPRISE SUPPORT	E-mail
	LOG IN

2. 在Dashboard页面,单击License Servers区域下CREATE LICENSE SERVER按钮。

۱.		$\sim$	~
٠.	1	-	
•		~	

÷	O é ullicensing.nvid	a.com	ė 🖈 🕭 🕹
0	NVIDIA LICENSING   D	lashboard	MMDM Aselection Hab   1   1   1   1   1   1   1   1   1
ដ	DASHBOARD	Entitlements MANAGE ENTITLEMENT	License Servers CREATE SERVERS CREATE SERVER
	ENTITLEMENTS	> ENTITLEMENT PAKID) / FEATURE ① START DATE EXPRATION ALLOCATED / TOTAL	> LICENSE SERVER BOUND SERVICE INSTANCE SERVICE INSTANCE TYPE
	LICENSE SERVERS	> management approved provide	You do not have any license servers, would you like to create one?
æ	SOFTWARE DOWNLOADS		U CREATE LICENSE SERVER
D	VIRTUAL GROUPS		
8	SERVICE INSTANCES		
10	LEASES	>	
ਿ	EVENTS	>	
83	USER MANAGEMENT	>	
0	MAINTENANCE		
	ENTERPRISE SUPPORT	>	
		>	
		>	

Entitlements区域下,是您目前拥有的全部License。

3. 在打开的Create License Server页面,填写相关信息,然后单击CREATE LICENSE SERVER。

Server Name	Feature	Lice	enses
Name this license server	Select a feature	✓ 1	
Description	Added Features		
Provide a short description	FEATURE		
	No featur	es have been added ye	t
MAC Address			
MAC Address (00:00:00:00:00:00:00 or XX-XX-XX-XX-XX-XX)			
Failover server configuration is optional. If configuring, you must provide a name AND MAC address			
Failover License Server			
Failover License Server Failover License Server			
Fallover License Server Fallover License Server Fallover MAC Address			

#### 必填参数说明如下:

参数	说明
Server Name	自定义您需要的License Server名称。
MAC Address	填写您用于搭建License Server的ECS实例的MAC地址。 您可以登录实例,使用 ipconfig -a 命令进行查询。
Feature	根据需要选择,并输入待添加的License数量,单击 <b>ADD</b> 。

4. 创建完成后,在License Servers页面单击 \_ 按钮,下载License文件。

Li	cense Server	^S	MANAGE SERVERS CREATE SERVER
≈	LICENSE SERVER	BOUND SERVICE INSTANCE	SERVICE INSTANCE TYPE
~	test LEGACY	n/a	n/a
		⊻ ∕ ⊝ 🤄	
	FEATURE		TOTAL
	NVIDIA RTX Virtual V (inkoys	Workstation-5.0	10

步骤五:导入License文件

1. 通过本地Web浏览器访问http://localhost:8080/licserver, 进入License Server管理界面。

⑦ 说明 请将访问地址中的*localhost* 替换为您的ECS实例的公网地址。



- 2. 在左侧导航栏的License Server区域,单击License Management。
- 3. 在License Management页面,单击Upload license file (.bin file)后的选择文件,选择本地的 License文件,然后单击打开(O)。

<mark> NVIDIA</mark> .			
License Management			
Logrand Conta     Browse for the loanse file you received from the IMDIA licensing portal, and then click Upload to     Description     Logrand Entropy Lipspe	rocess the Scense file.		
Log	• ля		×
Chemise Devel Manager	100 - 100	♥ / 11-286 / Pac / ○ 0 ■2214来	
> Estinos	1 24	▲ 名称 → 今天 (1)	传政日期 (*
	> 二 単一 出地部 > 30 対象	□ license_00 · · · · · · · · · · · · · · · · · ·	2022/1/27 15:48
	> 圖 税成 > ■ 週片 > ■ 元編	Contraction of the local division of the loc	
	> ● 下版 > ● 下版 > ♪ 音乐	There are	
	> <b>三</b> 皮耳 、 - ******		*
		文)华名(N)	所有文件 ~ 打开(0) 和3時

4. 单击Upload。

页面提示如下信息,表示License文件已导入成功。

License Management	
Successfully applied license file to license server.	
Brews for the License file you received from the NMDIA licensing portal, and then click Upload to process the license file.	
Uplaad license file ( Jain file):     [ 通信文件 ] 未进信任何文件	
	Cancel Upload

您可以在左侧导航栏的License Server区域,单击License Feature Usage,查看License数量以及使用情况。

			Search (cas	e-assative): Search by: Feature V 🔍 🔐
ated list. aga detaila.				
Version	Count	Available	Expiny	Wendor String
5.0	1	5		Quere and a second second
	ated list. Jage dekka. Version S 0	ane but. Nacionis <b>Court</b> 27 T	ant tic. ago antico. 12 Total Audulas	Servic San age India, Service Cause Andréa Conto 12 I I I I I I I I I I I I I I I I I I I

## 步骤六:测试网络连接和访问

以下操作以创建一台Windows操作系统的GPU虚拟化实例vgn6i为例,您也可以直接应用您已创建的GPU实例。

1. 创建一台GPU实例。

具体操作,请参见创建未配备驱动的GPU实例。

本示例选择的实例规格为GPU虚拟化实例vgn6i,操作系统为Windows Server 2019 数据中心版 64位中 文版。

2. 安装GRID驱动,然后添加License服务器并激活License。具体操作,请参见在GPU实例上安装GRID驱动 (Windows)。

激活License时,请填写您搭建了License Server的ECS实例公网IP,端口号为7070。

🛃 NVIDIA 控制面板	
文件(F) 编辑(E) 桌面(K) 帮助(H)	
🕝 fil 🔹 🕑 🍏	
选择一项任务	<b>停</b> 理许可证
□-3D 设置 通过预览调整图像设置	
□□- └─管理 3D 设置 □- 许可 □- └管理许可证	您可以通过应用许可证来启用附加功能。
	许可证类型: ✔ 您的系统已获 Quadro Virtual Data Center Workstation 许可。
	一级许可证服务器:
	7070
	二级许可证服务器:
	3mi
	说明: 主要许可证服务器聆听许可请求的端口编号。默认是 7070。

3. 前往License Server管理界面(*http://localhost:8080/licserver*),在左侧导航栏的License Server区域,单击Licensed Clients,即可查看到GPU实例信息,表示该NVIDIA vGPU软件客户端虚拟机目前正在使用License的功能。

Licensed Clients		
Licensed Clients with features consumed or reserved. Click	a Client ID for further details.	
Client ID	Client ID Type	Client Type
	E   Harrison	VIRTUAL
Go to page 1 v Total number of records, 1		
单击客户端ID,可	以查看该客户端的详细信息。	
Client Details		
Client Details		
Client ID: 0 Device ID Type: ETHERNET Client Type: VIRTUAL Client Expry: 2022		

# 2.在GPU实例上部署NGC环境

本文以搭建TensorFlow深度学习框架为例介绍如何在GPU实例上部署NGC环境。

### 前提条件

- 登录NGC网站,注册NGC账号。
- 登录NGC网站,获取NGC API key并保存到本地。登录NGC容器环境时需要验证您的NGC API Key。

#### 背景信息

NGC(NVIDIA GPU CLOUD)是NVIDIA开发的一套深度学习生态系统,可以使开发者免费访问深度学习软件堆 栈,建立适合深度学习的开发环境。

目前NGC在阿里云gn5实例作了全面部署,并且在镜像市场提供了针对NVIDIA Pascal GPU优化的NGC容器镜像。通过部署镜像市场的NGC容器镜像,开发者能简单快速地部署NGC容器环境,即时访问优化后的深度学习框架,极大程度缩减产品开发以及业务部署的时间,实现开发环境的预安装;同时支持调优后的算法框架,并且保持持续更新。

NGC网站提供了目前主流深度学习框架不同版本的镜像(例如Caffe、Caffe2、CNTK、MxNet、 TensorFlow、Theano、Torch),您可以选择需要的镜像部署环境。

支持部署NGC环境的实例规格族包括:

- gn4、gn5、gn5i、gn6v、gn6i、gn6e
- ebmgn5i、ebmgn6i、ebmgn6v、ebmgn6e

下面以gn5实例为例,为您演示创建GPU实例和部署NGC环境的步骤。

### 操作步骤

1. 创建一台gn5实例。具体操作,请参见使用向导创建实例。

在配置参数时,您需要注意以下几点:

- 地域:只能选择华北1(青岛)、华北2(北京)、华北3(张家口)、华北5(呼和浩特)、华东1(杭州)、华东2(上海)、华南1(深圳)。
- **实例**:选择gn5实例规格。
- 镜像: 单击镜像市场, 在弹出的对话框中, 找到NVIDIA GPU Cloud VM Image, 然后单击使用。

镜像市场[华北1]		×
	Q nvidia gpu cloud	授家
精选镜像	▲ 全部操作系统 ✓ 全部架构 ✓	
镜像分类 ∧ ✓ 全部 运行环境	NVIDIA GPU Cloud VM Image 基础系统:linux 架构:64位 NVIDIA GPU Cloud VM Image (虚拟机请像) 是运行针对NVIDIA ⑦	★★★★★ ¥0.00/月 13人已使用 使用
管理与监控		

○ 公网带宽:选择分配公网ⅠP地址。

⑦ 说明 如果这里没有分配公网IP地址,则在实例创建成功后,需要绑定EIP地址。

○ 安全组:选择一个安全组。安全组里必须开放TCP 22端口。如果您的实例需要支持HTTPS或DIGIT 6

服务,必须开放TCP 443 (用于HTTPS)或TCP 5000 (用于DIGITS 6)端口。

ECS实例创建成功后,请登录ECS管理控制台,记录实例的公网IP地址。

2. 连接ECS实例。

根据创建实例时选择的登录凭证选择以下任一方式连接ECS实例:

- 使用密码验证连接ECS实例
- 使用SSH密钥对验证连接ECS实例
- 3. 按界面提示输入NGC官网获取的NGC APIKey后按回车键,即可登录NGC容器环境。

? MobaXterm 8.4 ? (SSH client, X-server and networking tools) ➤ SSH session to ? SSH compression : -? SSH-browser ? X11-forwarding : < (remote display is forwarded through SSH) ? DISPLAY (automatically set on remote server) : 1 For more info, ctrl+click on <u>help</u> or visit our <u>website</u> Welcome to Ubuntu 16.04.4 LTS (GNU/Linux 4.4.0-116-generic x86\_64) \* Documentation: https://help.ubuntu.com https://landscape.canonical.com Management: Support: https://ubuntu.com/advantage Welcome to the NVIDIA GPU Cloud Virtual Machine. This environment is provided to enable you to easily run the Deep Learning containers from the NGC Registry. All of the documentation for how to use NGC and this VM are found at http://docs.nvidia.com/deeplearning/ngc Welcome to Alibaba Cloud Elastic Compute Service ! /usr/bin/xauth: file /root/.Xauthority does not exist lease enter your NGC APIkey to login to the NGC Registry:

4. 运行 nvidia-smi 命令。

您能查看当前GPU的信息,包括GPU型号、驱动版本等,如下图所示。

root@# nvidia-smi Thu Mar 29 20:50:01 2018	
NVIDIA-SMI 384.111 Driver Version: 384.111	
GPU Name Persistence-M Bus-Id Disp.A Volatile Fan Temp Perf Pwr:Usage/Cap  Memory-Usage GPU-Util	Uncorr. ECC   Compute M.
0 Tesla P100-PCIE 0ff   00000000:00:08.0 0ff   N/A 29C P0 27W / 250W   0MiB / 16276MiB   0%	0 Default
Processes: GPU PID Type Process name	GPU Memory   Usage
No running processes found	

- 5. 搭建TensorFlow深度学习框架。
  - i. 登录NGC网站,在TensorFlow镜像页面,获取 docker pull 命令。

Repositories	nvidia/tensorflow
nvidia 🧠	
caffe	docker pull nvcr.io/nvidia/tensorflow:18.03-py3
caffe2	
cntk	
cuda	
digits	
mxnet	
pytorch	
tensorflow	What is TensorFlow?
tensorrt	
theano	TensorFlow is an open source software library for numerical computation using data flow graphs. Nodes in the graph represent mathematical operations, while the graph edges represent the multidimensional
torch	data arrays (tensors) that flow between them. This flexible architecture lets you deploy computation to
hpc ^	one or more CPUs or GPUs in a desktop, server, or mobile device without rewriting code.

ii. 下载TensorFlow镜像。

docker pull nvcr.io/nvidia/tensorflow:18.03-py3

iii. 查看下载的镜像。

docker image ls

#### iv. 运行容器,完成TensorFlow开发环境的部署。

nvidia-docker run --rm -it nvcr.io/nvidia/tensorflow:18.03-py3



- 6. 选择以下任一种方式测试TensorFlow。
  - 简单测试TensorFlow。

#### python

```
import tensorflow as tf
hello = tf.constant('Hello, TensorFlow!')
sess = tf.Session()
sess.run(hello)
```

#### 如果TensorFlow正确加载了GPU设备,返回结果如下图所示。

```
root@^^^^ALL_ALL_BL///JNo
Python 3.5.2 (default, Nov 23 2017, 16:37:01)
[GCC 5.4.0 20160609] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import tensorflow as tf
>>> hello = tf.constant('Hello, TensorFlow!')
>>> sess = tf.Session()
2018-03-30 03:37:53.682157: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:892] s
be at least one NUMA node, so returning NUMA node zero
2018-03-30 03:37:53.682544: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1030] Foun
name: Tesla P100-PCIE-166B major: 0 memoryClockRate(GHz): 1.3285
pciBusID: 0000:00:08.0
totalMemory: 15.89GiB freeMemory: 15.60GiP
0010
   totalMemory: 15.89GiB freeMemory: 15.60GiB
2018-03-30 03:37:53.682583: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1120] Crea
16GB, pci bus id: 0000:00:08.0, compute capability: 6.0)
       >> sess.run(hello)
'Hello, TensorFlow!'
        >>
```

#### ○ 下载TensorFlow模型并测试TensorFlow。

git clone https://github.com/tensorflow/models.git cd models/tutorials/image/alexnet python alexnet\_benchmark.py --batch\_size 128 --num\_batches 100

#### 运行状态如下图所示。

conv1 [128, 56, 56, 64]
pool1 [128, 27, 27, 64]
conv2 [128, 27, 27, 192]
pool2 [128, 13, 13, 192]
conv3 [128, 13, 13, 384]
conv4 [128, 13, 13, 256]
conv5 [128, 13, 13, 256]
pool5 [128, 6, 6, 256]
2018-03-30 03:40:13.357785: I tensorflow/stream_executor/cuda/cuda_gpu_executor.cc:892] successful NUMA node read from SysFS
be at least one NUMA node, so returning NUMA node zero
2018-03-30 03:40:13.358207: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1030] Found device 0 with properties:
name: Tesla P100-PCIE-16GB major: 6 minor: 0 memoryClockRate(GHz): 1.3285
pciBusID: 0000:00:08.0
totalMemory: 15.89GiB freeMemory: 15.60GiB
2018-03-30 03:40:13.358245: I tensorflow/core/common_runtime/gpu/gpu_device.cc:1120] Creating TensorFlow device (/device:GPU:
16GB, pci bus id: 0000:00:08.0, compute capability: 6.0)
2018-03-30 03:40:15.916471: step 0, duration = 0.038
2018-03-30 03:40:16.299169: step 10, duration = 0.038
2018-03-30 03:40:16.682881: step 20, duration = 0.038
2018-03-30 03:40:17.065379: step 30, duration = 0.038
2018-03-30 03:40:17.448118: step 40, duration = 0.038
2018-03-30 03:40:17.830372: step 50, duration = 0.038
2018-03-30 03:40:18.213018: step 60, duration = 0.038
2018-03-30 03:40:18.595734: step 70, duration = 0.038
2018-03-30 03:40:18.978311: step 80, duration = 0.038
2018-03-30 03:40:19.361063: step 90, duration = 0.038
2018-03-30 03:40:19.705396: Forward across 100 steps, 0.038 +/- 0.000 sec / batch
2018-03-30 03:40:21.164735: step θ, duration = 0.090
2018-03-30 03:40:22.062778: step 10, duration = 0.090
2018-03-30 03:40:22.962202: step 20, duration = 0.090
2018-03-30 03:40:23.860856: step 30, duration = 0.090
2018-03-30 03:40:24.758891: step 40, duration = 0.090
2018-03-30 03:40:25.657170: step 50, duration = 0.090
2018-03-30 03:40:26.555194: step 60, duration = 0.090
2018-03-30 03:40:27.452843: step 70, duration = 0.090
2018-03-30 03:40:28.351092: step 80, duration = 0.090
2018-03-30 03:40:29.249606: step 90, duration = 0.090
2019 02 20 02, 10, 20 059090, Equipped basis and across 100 store 0,000 + ( 0,000 sec / batch

7. 保存TensorFlow镜像的修改。否则,下次登录时配置会丢失。

# 3.GPU AI模型训练最佳实践

本方案适用于AI图片训练场景,使用CPFS/NAS作为共享存储,利用容器服务Kubernetes版管理GPU云服务器集群进行图片AI训练。

## 实现的方式

- 搭建AI图片训练基础环境。
- 使用CPFS存储训练数据。
- 使用飞天AI加速训练服务加速训练。
- 使用Arena一键提交作业。

## 部署架构图



### 选用的产品

GPU服务器

GPU云服务器是基于GPU应用的计算服务,多适用于AI深度学习、视频处理、科学计算、图形可视化等应用场景。

更多关于GPU服务器的介绍,请参见GPU服务器产品详情页。

● 文件存储NAS

阿里云文件存储NAS是一个可共享访问、弹性扩展、高可靠、高性能的分布式文件系统。兼容POSIX文件 接口,可支持数千台计算节点共享访问,可以挂载到弹性计算ECS、弹性裸金属服务器、容器服务ACK、 弹性容器ECI、批量计算BCS、高性能计算EHPC、Al训练PAI等计算业务上提供高性能的共享存储,用户无 需修改应用程序,即可无缝迁移业务系统上云。

更多关于文件存储NAS的介绍,请参见文件存储NAS产品详情页。

● 文件存储CPFS

文件存储CPFS(Cloud Parallel File Storage),是阿里云完全托管、可扩展的并行文件存储系统,针对高性能计算场景的性能要求进行了深度优化,提供对数据毫秒级的访问和高聚合IO、高IOPS的数据读写请求,可以用于AI深度训练、自动驾驶、基因计算、EDA仿真、石油勘探、气象分析、机器学习、大数据分析以及影视渲染等业务场景中。

更多关于文件存储CPFS的介绍,请参见文件存储CPFS详情页。

#### ● 容器服务 ACK

容器服务Kubernetes版(ACK)提供高性能可伸缩的容器应用管理能力,支持企业级容器化应用的全生命 周期管理。整合阿里云虚拟化、存储、网络和安全能力,打造云端最佳容器化应用运行环境。

更多关于容器服务ACK的介绍,请参见容器服务 ACK产品详情页。

## 详细信息

点击查看最佳实践详情

## 更多最佳实践

点击查看更多阿里云最佳实践

# 4.在GPU实例上使用RAPIDS加速机器学 习任务

本文介绍了如何在GPU实例上基于NGC环境使用RAPIDS加速库,加速数据科学和机器学习任务,提高计算资源的使用效率。

## 背景信息

RAPIDS, 全称Real-time Acceleration Platform for Integrated Data Science, 是NVIDIA针对数据科学和机器学习推出的GPU加速库。更多RAPIDS信息,请参见官方网站。

NGC, 全称NVIDIA GPU CLOUD, 是NVIDIA推出的一套深度学习生态系统,供开发者免费访问深度学习和机器 学习软件堆栈,快速搭建相应的开发环境。NGC网站提供了RAPIDS的Docker镜像,预装了相关的开发环境。

JupyterLab是一套交互式的开发环境,帮助您高效地浏览、编辑和执行服务器上的代码文件。

Dask是一款轻量级大数据框架,可以提升并行计算效率。

本文提供了一套基于NVIDIA的RAPIDS Demo代码及数据集修改的示例代码,演示了在GPU实例上使用RAPIDS 加速一个从ETL到ML Training端到端任务的过程。其中,ETL时使用RAPIDS的cuDF,ML Training时使用 XGBoost。本文示例代码基于轻量级大数据框架Dask运行,为一套单机运行的代码。

② 说明 NVIDIA官方RAPIDS Demo代码请参见Mortgage Demo。

RAPIDS预装镜像已经发布到阿里云镜像市场,创建GPU实例时,您可以在镜像市场中搜索*NVIDIA RAPIDS*并 使用RAPIDS预装镜像。

⑦ 说明 该RAPIDS预装镜像使用Ubuntu 16.04 64-bit操作系统。

### 操作步骤

如果您创建GPU实例时使用了RAPIDS预装镜像,只需运行RAPIDS Demo,从启动JupyterLab服务开始操作即可。详情请参见启动JupyterLab服务。

如果您创建GPU实例时没有使用RAPIDS预装镜像,按照以下步骤使用RAPIDS加速机器学习任务:

- 步骤一: 获取NGC API Key
- 步骤二: 获取RAPIDS镜像下载命令
- 步骤三: 部署RAPIDS环境
- 步骤四:运行RAPIDS Demo

### 步骤一:获取NGC API Key

- 1. 在NGC注册页面注册NGC账号。
- 2. 登录NGC网站。
- 3. 单击页面右上角用户名,然后单击Setup。
- 4. 在Setup页面的Generate API Key区域,单击Get API Key。
- 5. 在API Key页面, 单击Generate API Key。
- 6. 在Generate a New API Key对话框中, 单击Confirm。

② 说明 新的NGC API Key会覆盖旧的NGC API Key。如果您已持有NGC API Key,请确保不再需要 旧的NGC API Key。

7. 复制APIKey并保存到本地。



# 步骤二:获取RAPIDS镜像下载命令

完成以下操作,获取RAPIDS镜像下载命令:

- 1. 登录NGC网站。
- 2. 在页面左侧导航栏,选择CATALOG > Explore Catalog。
- 3. 在NVIDIA NGC: AI Development Catalog页面的搜索栏中,输入RAPIDS。
- 4. 在搜索结果中,单击RAPIDS镜像。

Y Query: RAPIDS X				X Sort: Relevance V		
DVIDIA. DEEP LEARNING INSTITUTE	RAPIDS	Accelerated with	RAPIDS	Accelerated with	RAPIDS	Accelerated with
OLI RAPIDS Course - Base Environment	RAPIDS Container		RAPIDS Cloud Machi Container	ne Learning	GPU Accelerated ML Container	workflows with R
Base environment used in the NVIDIA Deep Learning Institute (DLI) Course Fundamentals of Accelerated Data Science with RAPIDS, along with Next	The RAPIDS suite of software libraries gives you the freedom to execute end-to- end data science and analytics pipelines entirely on GPUs.		RAPIDS is a suite of o that bring GPU accel science pipelines. Us based machine learn	pen-source libraries eration to data ers building cloud- ing experiments ca	Demonstration of GP Machine Learning Da workflows using RAPI	U Accelerated ta Science DS.
liew Labels Pull Tag	View Labels Pull T	ag	View Labels Pull T	ag	View Labels Pull Ta	9

5. 获取docker pull命令。

本文示例代码基于RAPIDS 0.8版本镜像编写,因此在运行本示例代码时,使用Tag为0.8版本的镜像。实际操作时,请选择您匹配的版本。

i. 单击Tags页签。

Catalog: Containers / Containers: nvidia:rapidsai:rapidsai										
RAPIDS										
Publisher	Built By	Latest Tag	Modified	Size						
Open Source	NVIDIA	cuda10.1-base	December 11,	3.97 GB						
Multinode Support	Multi-Arch Support									
NO										
Description										
The RAPIDS suite of se entirely on GPUs.	oftware libraries gives you t	he freedom to execute er	nd-to-end data science and	d analytics pipelines						
Labels										
Covid-19 Machine I	earning									
Pull Command	Pull Command									
	o/nvidia/rapidsai/rapids.	ai:cuda10.1-base-centos		D						
Overview Tag	s Layers Relat	ed Collections								

ii. 找到并复制Tag信息。本示例中,选择 0.8-cuda10.0-runtime-ubuntu16.04-gcc5-py3.6 。

0.8-cuda10.0-runtime-ubuntu16.04-gcc5-py3.6 07/11/2019 07:10AM		2.96 GB	⊥ Pull Tag
Manifest			
DIGEST	OS/ARCH	COMPRESSED SIZE	CREATED
sha256:9b0e30899485faf554b8a8a9a7465e57abab000bef9fbc0d	linux/amd64	2.96 GB	07/11/2019

 iii. 返回页面顶部,复制Pull Command中的命令到文本编辑器,将镜像版本替换为对应的Tag信息, 并保存。

本示例中,将 cuda9.2-runtime-ubuntu16.04 替换为 0.8-cuda10.0-runtime-ubuntu16.04-gcc 5-py3.6 。

保存的docker pull命令用于下载RAPIDS镜像。关于如何下载RAPIDS镜像,请参见步骤三:部署 RAPIDS环境。

	Publisher	Built By	Latest Tag	Modified	Size
	Open Source	NVIDIA	cuda9.2-runtime	July 11, 2019	2.94 GB
	Description				
	The RAPIDS suite of soft	vare libraries gives you the	freedom to execute end-to-end	data science and analytic	s pipelines entirely on GPUs.
	Labels				
RAFIDS	Machine Learning				
i i i i i i i i i i i i i i i i i i i	Pull Command				
	docker pull nvcr.io/	2 widia/rapidsai/rapidsai:	tuda9.2-runtime-ubuntul6.04		•••

### 步骤三: 部署RAPIDS环境

完成以下操作,部署RAPIDS环境:

1. 创建一台GPU实例。详细步骤请参见使用向导创建实例。

参数配置说明如下:

- 实例: RAPIDS仅适用于特定的GPU型号(采用NVIDIA Pascal及以上架构),因此您需要选择GPU型号符合要求的实例规格,目前有gn6i、gn6v、gn5和gn5i,详细的GPU型号请参见实例规格族。建议您选择显存更大的gn6i、gn6v或gn5实例。本示例中,选用了显存为16 GB的GPU实例。
- 镜像:在镜像市场中搜索并使用 NVIDIA GPU Cloud VM Image 。

境像市场[华东 2 (上海)] ×								
			R, NVIDIA GPU Cloud		搜索			
	精选镜像	*	全部操作系统 🗸 全部架构 🗸					
	镜像分类 // ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓	^	NVIDIA GPU Cloud VM Image 基础系统:linux 架构:64位 NVIDIA GPU Cloud VM Image (虚拟机确像) 是运行针对NVIDIA	9.07.0 🔻 ⊘	★★★★★ 1974人已使用	¥ 0.00/月 使用		
	运行环境 管理与监控							
	建站系统							

- 公网带宽:选择分配公网IPv4地址或者在实例创建成功后绑定EIP地址。具体操作,请参见弹性公网IP文档中的绑定ECS实例。
- 安全组:选择的安全组需要开放以下端口:
  - TCP 22 端口,用于SSH登录
  - TCP 8888端口,用于支持访问JupyterLab服务
  - TCP 8787端口、TCP 8786端口,用于支持访问Dask服务

- 2. 连接GPU实例。连接方式请参见连接方式介绍。
- 3. 输入NGC APIKey后按回车键,登录NGC容器环境。



4. (可选)运行nvidia-smi查看GPU型号、GPU驱动版本等GPU信息。

建议您了解GPU信息,预判规避潜在问题。例如,如果NGC的驱动版本太低,新Docker镜像版本可能会 不支持。

5. 运行docker pull命令下载RAPIDS镜像。

关于如何获取docker pull命令,请参见步骤二:获取RAPIDS镜像下载命令。

```
docker pull nvcr.io/nvidia/rapidsai/rapidsai:0.8-cuda10.0-runtime-ubuntu16.04-gcc5-py3.
```

6. (可选)查看下载的镜像。

建议您查看Docker镜像信息,确保下载了正确的镜像。

docker images

7. 运行容器部署RAPIDS环境。

```
docker run --runtime=nvidia \
    --rm -it \
    -p 8888:8888 \
    -p 8787:8787 \
    -p 8786:8786 \
    nvcr.io/nvidia/rapidsai/rapidsai:0.8-cuda10.0-runtime-ubuntu16.04-gcc5-py3.6
```

### 步骤四:运行RAPIDS Demo

完成以下操作,运行RAPIDS Demo:

1. 在GPU实例上下载数据集和Demo文件。

```
# Get apt source address and download demos.
source_address=$(curl http://100.100.200/latest/meta-data/source-address|head -n 1)
source_address="${source_address}/opsx/ecs/linux/binary/machine_learning/"
cd /rapids
wget $source_address/rapids_notebooks_v0.8.tar.gz
tar -xzvf rapids_notebooks_v0.8.tar.gz
cd /rapids/rapids_notebooks_v0.8/xgboost
wget $source_address/data/mortgage/mortgage_2000_lgb.tgz
```

#### 2. 在GPU实例上启动JupyterLab服务。

#### 推荐直接使用命令启动。

```
# Run the following command to start JupyterLab and set the password.
cd /rapids/rapids_notebooks_v0.8/xgboost
jupyter-lab --allow-root --ip=0.0.0.0 --no-browser --NotebookApp.token='YOUR PASSWORD'
# Exit JupyterLab.
sh ../utils/stop-jupyter.sh
```

- 除使用命令外,您也可以执行脚本 sh ../utils/start-jupyter.sh 启动jupyter-lab,此时无法设 置登录密码。
- 您也可以连续按两次 Ctrl+C 退出JupyterLab服务。
- 3. 打开浏览器,在地址栏输入 http://*您的GPU实例IP地址*:8888 远程访问JupyterLab。
  - ⑦ 说明 推荐使用Chrome浏览器。

如果您在启动JupyterLab服务时设置了登录密码, 会跳转到密码输入界面。

	💭 Jupyte	er	
Password or token:	1	Log in	
Token authentication is enab	led		
If no password has been configured, yo the URL, or paste it above. This require	u need to open th ment will be lifted	e notebook server with its login to if you <u>enable a password</u> .	ken i

4. 运行NoteBook代码。

该案例是一个抵押贷款回归的任务,更多信息,请参见代码执行过程。登录成功后,可以看到 NoteBook代码包括以下内容:

• *xgboost\_E2E.ipynb*文件: XGBoost Demo文件。双击文件可以查看文件详情,单击下图中的执行按 钮可以逐步执行代码,每次执行一个Cell。

$\Box$	File Edit View	Run I	Kernel Tal	os Settings Help								
1.00	+	10	±	C	<b>13</b> (	aunc	her			:	< 1	■ xgboost_E2E.ipynb ×
_	<b>. +</b>			8	+	Ж	Ō	Ê	► 1		C Markdown v	
ġ.	Name		•	Last Modified				c		(0) 20	0.0	010 Alikaka Casur Haldina Limitad
	• 🗖 xgboost_E2E.ip	ynb		6 days ago	Copyright (C) 2010-2019 Alibaba Group Holding Limited					DIS Alibaba Group Holding Limited		
۲	🗅 mortgage_200	0_1gb.tgz		2 months ago				Сору	right	(C) 201	8 N\	/IDIA Corporation
	README.md			6 days ago	1.1							
2								M	ort	dad	ae	Workflow
											,-	
								The	e D	atas	et	
								The c	latas	et used	with	n this workflow is derived from Fannie Mae's Single-Family Loan

○ *mortgage\_2000\_1gb.tgz*文件: 2000年的抵押贷款回归训练数据(1G分割的perf文件夹下的文件不 会大于1 G,使用1 G分割的数据可以更有效的利用GPU显存)。

## 代码执行过程

该案例基于XGBoost演示了数据预处理到训练的端到端的过程,主要分为三个阶段:

- ETL(Extract-Transform-Load): 主要在GPU实例上进行。将业务系统的数据经过抽取、清洗转换之后 加载到数据仓库。
- Data Conversion:在GPU实例上进行。将在ETL阶段处理过的数据转换为用于XGBoost训练的DMatrix格式。
- ML-Training: 默认在GPU实例上进行。使用XGBoost训练梯度提升决策树。

NoteBook代码的执行过程如下:

1. 准备数据集。

本案例的Shell脚本会默认下载2000年的抵押贷款回归训练数据(mortgage\_2000\_1gb.tgz)。

如果您想获取更多数据用于XGBoost模型训练,可以设定参数download\_url指定下载路径,具体下载地 址请参见Mortgage Data。

### 示例效果如下:



#### 2. 设定相关参数。

参数名称	说明
start_year	指定选择训练数据的起始时间,ETL时会处理start_year到end_year之间的数据。
end_year	指定选择训练数据的结束时间,ETL时会处理start_year到end_year之间的数据。
train_with_gpu	是否使用GPU进行XGBoost模型训练,默认为 <i>True</i> 。
gpu_count	指定启动worker的数量,默认为1。您可以按需要设定参数值,但不能超出GPU实例 的GPU数量。
part_count	指定用于模型训练的performance文件的数量,默认为 <i>2 * gpu_count</i> 。如果参数值 过大,在Dat <i>a</i> Conversion阶段会报错超出GPU内存限制,错误信息会在NoteBook 后台输出。

#### 示例效果如下:

[5];	acq_dsta_path = "()/acq", format(mortgage_dir) perf_dsta_path = '()/perf_'format(mortgage_dir) col_mems_path = '()/mesc.cv', format(mortgage_dir)
	<pre>start_year = 2000 end_year is inclusive</pre>
	# whether use GPUs for XGBoost training train_with.gpu = True
	# The member of GPUs to be used, value range: I to get_gpu_mas(). Default value: 1. # # This parameter would use for starting dask-worker, doing ETL, doing Conversion and training model(if train_with_gpu=True). gpu_count 1 i
	<pre># The number of performance files in the perf folder part_number = len(os.listdir(perf_data_path))</pre>
	# if your download 100 Split troin data(the filename and with '1gb.gr'), each performance file is no large than 100, # in this comple, a 100 GPU can process 2 or 3 preformance files. By default, one GPU is set to process 2 files. part_compl = 2 * gpu_count if a performance files. The gpu_count is a performance file.
	print(')>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>
	>>> Using "1" GPU(GPUs). >>> Ell - process performance files from "2000" to "2000". >>>> Data Conversion - select "2" performance data processed in the EIL phase and convert to DMutrix-format for XGBoost training. >>>> N L - Whether to use the GMP confidence files "True".

#### 3. 启动Dask服务。

代码会启动Dask Scheduler,并根据gpu\_count参数启动worker用于ETL和模型训练。启动Dask服务 后,您也可以通过Dask Dashboard直观地监控任务,打开方法请参见Dask Dashboard。

#### 示例效果如下:

<pre># run dask-worker cmd = "hostmameall-ip-addresses" process = subprocess.Popen(cmd.split output, error = process.communicate IPADDR = striculptut.dccdd().split cluster = localCUDACluster(n_workers client = Client(cluster) client</pre>	<pre>(), stdout=subprocess.PIPE) ) [0] agpu_count, ip=IPADOR)</pre>
Client Client Client Client Dashboard: http://172.17.0.243894	Cluster Workers: 1 Cores: 1 • Memory: 507:25 G8

#### 4. 启动ETL。

ETL阶段会进行到表关联、分组、聚合、切片等操作,数据格式采用cuDF库的DataFrame格式(类似于 pandas的DataFrame格式)。

#### 示例效果如下:

ETL
Perform all of ETL with a single call to
<pre>process_quarter_gpu(year=year, quarter=quarter, perf_file=file)</pre>
%%time
# NOTE: The EIL calculates additional features which are then dropped before creating the XGBoost DMatrix. # This can be optimized to avoid calculating the dropped features.
<pre>gpu_dfs = [] gpu_time = 0</pre>
quarter = 1
count = 0
while year <= end year:
<pre>for file in glob(os.path.join(perf_data_path + "/Performance_" + str(year) + "Q" + str(quarter) + """)):     gpu_dfs.append(process_quarter_gpu(yearsyear, quartersquarter, perf_filesfile))     count += 1</pre>
quarter += 1
if quarter == 5:
year += 1
quarter = 1
wait(gpu_dfs)
CPU times: user 560 ms, sys: 28 ms, total: 588 ms
Wall time: 20.9 s

#### 5. 启动Data Conversion。

将DataFrame格式的数据转换为用于XGBoost训练的DMatrix格式,每个worker处理一个DMatrix对象。

#### 示例效果如下:

Load the data from host memory, and convert to CSR
%%time
<pre>gpu_dfs = [delayed(DatsFrame.from_arrow)(gpu_df) for gpu_df in gpu_dfs[:part_count]] gpu_dfs = [gpu_df for gpu_df in gpu_dfs] wait(gpu_dfs)</pre>
<pre>tmp_map = [(gpu_df, list(client.who_has(gpu_df).values())[0]) for gpu_df in gpu_dfs] new map = {}</pre>
for key, value in tmp_map:
if value not in new_map:
new_map[value] = [key]
else:
new_map[value].append(key)
del(tmp_map)
gpu_dfs = []
for list_delayed in new_map.values():
<pre>gpu_dfs.append(delayed(cudf.concat)(list_delayed))</pre>
del(new map)
<pre>gpu_dfs = [(gpu_df[['delinquency_12']], gpu_df[delayed(list)(gpu_df.columns.difference(['delinquency_12']))]) for gpu_df in gpu_dfs] gpu_dfs = [(gpu_df[0].persist(), gpu_df[1].persist()) for gpu_df in gpu_dfs]</pre>
enu dfs = [dask_delaved(xeb_DMatrix)(enu df[1], enu df[2]) for enu df in enu dfs]
goud fs = [goud f.persist() for goud f in goud fs
gc.collect()
wait(gpu_dfs)
CPU times: user 200 ms, sys: 4 ms, total: 204 ms
Wall time: 4.3 s

#### 6. 启动MLTraining。

使用dask-xgboost启动模型训练,dask-xgboost负责多个dask worker间的通信协同工作,底层仍然调用xgboost执行模型训练。

#### 示例效果如下:



# Dask Dashboard

Dask Dashboard支持任务进度跟踪、任务性能问题识别和故障调试。

Dask服务启动后,在浏览器地址栏中访问 http://*您的GPU实例IP地址*:8787/status 即可进入Dashboard主 界面。



# 相关函数

函数功能	函数名称
下载文件	def download_file_from_url(url, filename):
解压文件	def decompress_file(filename, path):
获取当前机器的GPU个数	def get_gpu_nums():
管理GPU内存	<ul> <li>def initialize_rmm_pool():</li> <li>def initialize_rmm_no_pool():</li> <li>def run_dask_task(func, **kwargs):</li> </ul>
提交DASK任务	<ul> <li>def process_quarter_gpu(year=2000, quarter=1, perf_file=""):</li> <li>def run_gpu_workflow(quarter=1, year=2000, perf_file="", **kwargs):</li> </ul>

函数功能	函数名称	
使用cuDF从CSV中加载数据	<ul> <li>def gpu_load_performance_csv(performance_path, **kwargs):</li> <li>def gpu_load_acquisition_csv(acquisition_path, **kwargs):</li> <li>def gpu_load_names(**kwargs):</li> </ul>	
处理和提取训练数据的特征	<ul> <li>def null_workaround(df, **kwargs):</li> <li>def create_ever_features(gdf, **kwargs):</li> <li>def join_ever_delinq_features(everdf_tmp, delinq_merge, **kwargs):</li> <li>def create_joined_df(gdf, everdf, **kwargs):</li> <li>def create_12_mon_features(joined_df, **kwargs):</li> <li>def combine_joined_12_mon(joined_df, testdf, **kwargs):</li> <li>def final_performance_delinquency(gdf, joined_df, **kwargs):</li> <li>def join_perf_acq_gdfs(perf, acq, **kwargs):</li> <li>def last_mile_cleaning(df, **kwargs):</li> </ul>	

# 相关文档

• 在GPU实例上使用RAPIDS加速图像搜索任务

# 5.在GPU实例上使用RAPIDS加速图像搜 索任务

本文以使用RAPIDS加速图像搜索任务为例,介绍如何在预装镜像的GPU实例上使用RAPIDS加速库。

### 前提条件

使用本教程进行操作前,请确保您已经注册了阿里云账号。如还未注册,请先完成账号注册。

## 背景信息

RAPIDS, 全称Real-time Acceleration Platform for Integrated Data Science, 是NVIDIA针对数据科学和机器学习推出的GPU加速库。更多RAPIDS信息请参见官方网站。

基于图像识别和搜索,图像搜索任务可以实现以图搜图,在不同行业应用和业务场景中帮助您搜索相同或相 似的图片。

图像搜索任务背后的两项主要技术是特征提取及向量化、向量索引和检索。本文案例中,使用开源框架 TensorFlow和Keras配置生产环境,然后使用ResNet50卷积神经网络完成图像的特征提取及向量化,最后使 用RAPIDS cuML库的KNN算法实现BF方式的向量索引和检索。

⑦ 说明 BF (Brute Force)检索方法是一种百分百准确的方法,对距离衡量算法不敏感,适用于所有的距离算法。

本文案例在阿里云gn6v(NVIDIA Tesla V100)实例上执行。执行案例后,对比了GPU加速的RAPIDS cuml KNN与CPU实现的scikit-learn KNN的性能,可以看到GPU加速的KNN向量检索速度为CPU的近600倍。

本文案例为单机单卡的版本,即一台GPU实例搭载一块GPU卡。

### 操作步骤

执行以下操作完成一次图像搜索任务:

- 1. 创建GPU实例
- 2. 启动和登录JupyterLab
- 3. 执行图像搜索案例

### 步骤一: 创建GPU实例

具体步骤请参见使用向导创建实例。

- **实例**: RAPIDS仅适用于特定的GPU型号(采用NVIDIA Pascal及以上架构),因此您需要选择GPU型号符合 要求的实例规格,目前有gn6i、gn6v、gn5和gn5i。本文案例中,选用了ecs.gn6v-c8g1.2xlarge实例规 格。
- 镜像:在镜像市场中使用关键字 RAPIDS,搜索并使用预装了 RAPIDS加速库的镜像。

镜像市场[华东1(桥	ðH)]	×
	۹. RAPIDS	
精选镜像	▲ 全部操作系统 ∨ 全部架构 ∨	
嶺像分类 ✓ 全部 操作系统	<ul> <li>◇ Ubuntu16.04(预装NVIDIA RAPIDS)</li> <li>基础系统: linux 架构: 64位</li> <li>该请像使用Ubuntu16.04 64bit系统,预装NVIDIA RAPIDS机器学</li> </ul>	¥0.00/时 使用
运行环境 管理与监控		
建站系统 应用开发		

• 安全组:选择的安全组需要开放TCP 8888端口,用于支持访问JupyterLab服务。

# 步骤二:启动和登录JupyterLab

1. 连接GPU实例,运行以下命令启动JupyterLab服务。

⑦ 说明 连接GPU实例的步骤请参见连接方式概述ECS远程连接操作指南。

```
# Go to the notebooks directory.
cd /rapids
# Run the following command to start JupyterLab and set the logon password:
jupyter-lab --allow-root --ip=0.0.0.0 --no-browser --NotebookApp.token='your logon pass
word'
# Exit jupyterlab: press Ctrl+C twice.
```

2. 在您的本地机器上打开浏览器。输入 http://(IP address of your GPU instance):8888 远程访问 JupyterLab。

? 说明 推荐使用Chrome浏览器。

3. 输入启动命令中设置的密码, 然后单击Log in。

💭 jupyt	er	
Password or token:	Log in	
Token authentication is enabled		
If no password has been configured, you need to open to the URL, or paste it above. This requirement will be lifted	he notebook server with its login toker d if you <u>enable a password</u> .	n in

## 步骤三:执行图像搜索案例

- 1. 进入案例所在目录 rapids\_notebooks\_v0.7/cuml。
- 2. 双击cuml\_knn.ipynb文件。
- 3. 单击 🕨 。

File E	dit View	Run Kernel	Tabs	Settings Help								
	+	10	<u>*</u>	C	🖾 Lau	ncher		- (	3) ×		cuml_knn.ipynb	×
1 <b>त</b> > ra	oids_notebo	oks_v0.7 > cuml			8	- *	Ō	0	•	C	Markdown ~	
Name			•	Last Modified								
🗖 dat	3			a day ago			Im	an	- 54	ar	ch	
2 🗖 cur	ıl_knn.ipynb	]		23 minutes ago				ag		a	ch	
♥ RE/	DME.md			3 days ago			Int	rod	uctio	on		
							The o	lemo i	s comp	osed	of three core phase	25:
							1.1	Datase	t			
							2.1	eature	2			
							3.1	KNN				

# 案例执行过程

图像搜索案例的执行过程分为三个步骤:处理数据集、提取图片特征和搜索相似图片。本文案例结果中对比了GPU加速的RAPIDS cuml KNN与CPU实现的scikit-learn KNN的性能。

- 1. 处理数据集。
  - i. 下载和解压数据集。

本文案例中使用了STL-10数据集,该数据集中包含10万张未打标的图片,图片的尺寸均为:96 x 96 x 3。您可以使用其他数据集,为便于提取图片特征,请确保数据集中图片的尺寸相同。

本文案例提供了 download\_and\_extract(data\_dir) 方法供您下载和解压STL-10数据集。RAPIDS 镜像中已经将数据集下载到./data目录,您可以执行 download\_and\_extract() 方法直接解压数 据集。



ii. 读取图片。

从数据集解压出的数据为二进制格式,执行 read\_all\_images(path\_to\_data) 方法加载数据并转 换为NHWC(batch, height, width, channels)格式,以便用Tensorflow提取图片特征。

Read Data

```
[3]: # the path of unlabeled data
path_unlabeled = os.path.join(data_dir, 'stl10_binary/unlabeled_X.bin')
# get images from binary
images = read_all_images(path_unlabeled)
print('>>> images shape: ', images.shape)
>>> images shape: (100000, 96, 96, 3)
```

#### iii. 展示图片。

执行 show image (image) 方法随机展示一张数据集中的图片。

	Show Image
[4]:	<pre>import random import matplotlib.pyplot as plt %matplotlib inline</pre>
	<pre>def show_image(image):     """show image"""     fig = plt_figure(figsize=(3, 3))</pre>
	<pre>plt.imshow(image) plt.show() fig.clear()</pre>
[10]:	<pre># random show a image rand_image_index = random.randint(0, images.shape[0]) show_image(images[rand_image_index])</pre>
	80 - 20 40 60 80

iv. 分割数据集。

按照9:1的比例把数据集分为两部分,分别用于创建图片索引库和搜索图片。



2. 提取图片特征。

使用开源框架Tensorflow和Keras提取图片特征,其中模型为基于ImageNet数据集的 ResNet50(notop)预训练模型。

#### i. 设定Tensorflow参数。

Tensorflow默认使用所有GPU显存,我们需要留出部分GPU显存供cuML使用。您可以选择一种方法 设置GPU显存参数:

■ 方法1: 依据运行需求进行显存分配。

config.gpu\_options.allow\_growth = True

■ 方法2: 设定可以使用的GPU显存比例。本案例中使用方法2,并且GPU显存比例默认设置为0.3, 即Tensorflow可以使用整块GPU显存的30%,您可以依据应用场景修改比例。

config.gpu\_options.per\_process\_gpu\_memory\_fraction = 0.3

Image Features
<pre># set tensorflow params to adjust GPU memory usage, if use default params, tensorflow would use # nearly all of the gpu memory, we need reserve some gpu memory for cuml. import os # only use device 0 os.environ["CUDA_VISIBLE_DEVICES"] = "0"</pre>
<pre>import tensorflow as tf from keras.backend.tensorflow_backend import set_session config = tf.ConfigProto() # method 1: allocate gpu memory base on runtime allocations # config.gpu_options.allow_growth = True # method 2: determines the fraction of the onerall amount of memory # that each visibel GPU should be allocated. config.gpu_options.per_process_gpu_memory_fraction = 0.3 set_session(tf.Session(config=config))</pre>
Using TensorFlow backend.

### ii. 下载ResNet50 (notop) 预训练模型。

连接公网下载模型(大小约91 M),下载完成后默认保存到/root/.keras/models/目录。

参数名称	说明
weights	取值范围: <ul> <li>None:随机初始化权重值。</li> <li><i>imagenet</i>:权重值的初始值设置为通过ImageNet预训练过的模型的权重值。</li> </ul> 本案例中设置为 <i>imagenet</i> 。
include_top	取值范围: <ul> <li>True:包含整个ResNet50网络结构的最后一个全链接层。</li> <li>False:不包含整个ResNet50网络结构的最后一个全链接层。</li> <li>本案例中,使用神经网络模型ResNet50的主要目的是提取图片特征而非分类图片,因此设置为<i>False</i>。</li> </ul>
input_shape	可选参数,用于设置图片的输入shape,仅在include_top设置为 <i>False</i> 时生效。 您必须为图片设置3个inputs channels,且宽和高不应低于32。此处设为 <i>(96, 9 6, 3)</i> 。
pooling	在include_top设置为 <i>False</i> 时,您需要设置池化层模式,取值范围: <ul> <li><i>None</i>:输出为4D tensor。</li> <li><i>avg</i>:平均池化,输出为2D tensor。</li> <li><i>max</i>:最大池化,输出为2D tensor。</li> </ul> 本案例中设置为 <i>max</i> 。

### 您可以执行 model.summary() 方法查看模型的网络结构。

(None,	3, 3, 2048)	0	bn5c_branch2c[0][0] activation_46[0][0]
(None,	3, 3, 2048)	0	add_16[0][0]
(None,	2048)	0	activation_49[0][0]
	(None, (None, (None,	(None, 3, 3, 2048) (None, 3, 3, 2048) (None, 2048)	(None, 3, 3, 2048) 0 (None, 3, 3, 2048) 0 (None, 2048) 0

#### iii. 提取图片特征。

对分割得到的两个图片数据集执行 model.predict() 方法提取图片特征。

- 3. 搜索相似图片。
  - i. 使用cuml KNN搜索相似图片。

. . . . . . .

通过 k=3 设置K值为3,即查找最相似的3张图片,您可以依据使用场景自定义K值。

```
其中, knn_cuml.fit() 方法为创建索引阶段, knn_cuml.kneighbors() 为搜索近邻阶段。
```

	cumi KNN
[12]:	<pre>from cuml.neighbors import NearestNeighbors</pre>
[13]:	<pre>%%time knn_cuml = NearestNeighbors() knn_cuml.fit(train_features)</pre>
	CPU times: user 888 ms, sys: 60 ms, total: 948 ms Wall time: 192 ms
[14]:	<pre>%%time distances_cuml, indices_cuml = knn_cuml.kneighbors(query_features, k=3)</pre>
	CPU times: user 1.59 s, sys: 492 ms, total: 2.08 s Wall time: 791 ms

KNN向量检索耗时791 ms。

ii. 使用scikit-learn KNN搜索相似图片。

通过 n neighbors=3 设置K值为3,通过 n jobs=-1 设置使用所有CPU进行近邻搜索。

⑦ 说明 ecs.gn6v-c8g1.2xlarge的配置为8 vCPU。

	sklearn KNN
[15]:	<pre>from sklearn.neighbors import NearestNeighbors</pre>
[16]:	<pre>%%time knn_sk = NearestNeighbors(n_neighbors=3, metric='sqeuclidean', n_jobs=-1) knn_sk.fit(train_features)</pre>
	CPU times: user 856 ms, sys: 36 ms, total: 892 ms Wall time: 114 ms
[17]:	<pre>%%time distances_sk, indices_sk = knn_sk.kneighbors(query_features, 3)</pre>
	CPU times: user 18.2 s, sys: 29.9 s, total: 48.1 s Wall time: 7min 34s

KNN向量检索耗时7分34秒。

iii. 对比cuml KNN和scikit-learn KNN的搜索结果。

对比两种方式的KNN向量检索速度,使用GPU加速的cuml KNN耗时791 ms,使用CPU的scikit-learn KNN耗时7min 34s。前者为后者的近600倍。

验证两种方式的输出结果是否相同,输出结果为两个数组:

- distance: 最小的K个距离值。本案例中搜索了10000张图片,K值为3,因此 distance.shape=(
  10000,3) 。
- indices: 对应的图片索引。 indices.shape=(10000, 3) 。

由于本案例所用数据集中存在重复图片,容易出现图片相同但索引不同的情况,因此使用distances,不使用indices对比结果。考虑到计算误差,如果两种方法得出的10000张图片中的3 个最小距离值误差都小于1,则认为结果相同。



## 图片搜索结果

本案例从1万张搜索图片中随机选择5张图片并搜索相似图片,最终展示出5行4列图片。

第一列为搜索图片,第二列至第四列为图片索引库中的相似图片,且相似性依次递减。每张相似图片的标题 为计算的距离,数值越大相似性越低。



# 6.RAPIDS加速机器学习最佳实践

本方案适用于使用RAPIDS加速库和GPU云服务器来对机器学习任务或者数据科学任务进行加速的场景。相比 CPU,利用GPU和RAPIDS在某些场景下可以取得非常明显的加速效果。

## 解决的问题

- 搭建RAPIDS加速机器学习环境。
- 使用容器服务Kubernetes版部署RAPIDS环境。
- 使用NAS存储计算数据。

## 部署架构图



## 选用的产品

● GPU服务器

GPU云服务器是基于GPU应用的计算服务,多适用于AI深度学习、视频处理、科学计算、图形可视化等应用场景。

更多关于GPU服务器的介绍,请参见GPU服务器产品详情页。

● 文件存储NAS

阿里云文件存储NAS是一个可共享访问、弹性扩展、高可靠、高性能的分布式文件系统。兼容POSIX文件 接口,可支持数千台计算节点共享访问,可以挂载到弹性计算ECS、弹性裸金属服务器、容器服务ACK、 弹性容器ECI、批量计算BCS、高性能计算EHPC、Al训练PAI等计算业务上提供高性能的共享存储,用户无 需修改应用程序,即可无缝迁移业务系统上云。

更多关于文件存储NAS的介绍,请参见文件存储NAS产品详情页。

● 容器服务 ACK

容器服务Kubernetes版(ACK)提供高性能可伸缩的容器应用管理能力,支持企业级容器化应用的全生命 周期管理。整合阿里云虚拟化、存储、网络和安全能力,打造云端最佳容器化应用运行环境。 更多关于容器服务ACK的介绍,请参见<mark>容器服务 ACK产品详情页</mark>。

# 详细信息

点击查看最佳实践详情

更多最佳实践

点击查看更多阿里云最佳实践

# 7.RAPIDS加速图像搜索最佳实践

本方案适用于使用RAPIDS加速平台和GPU云服务器来对图像搜索任务进行加速的场景。相比CPU,利用 GPU+RAPIDS在图像搜索场景下可以取得非常明显的加速效果。

## 解决的问题

- 搭建RAPIDS加速图像搜索环境。
- 使用容器服务Kubernetes版部署图像搜索环境。
- 使用NAS存储计算数据。

## 部署架构图



## 选用的产品

GPU服务器

GPU云服务器是基于GPU应用的计算服务,多适用于AI深度学习、视频处理、科学计算、图形可视化等应用场景。

更多关于GPU服务器的介绍,请参见GPU服务器产品详情页。

● 文件存储NAS

阿里云文件存储NAS是一个可共享访问、弹性扩展、高可靠、高性能的分布式文件系统。兼容POSIX文件 接口,可支持数千台计算节点共享访问,可以挂载到弹性计算ECS、弹性裸金属服务器、容器服务ACK、 弹性容器ECI、批量计算BCS、高性能计算EHPC、Al训练PAI等计算业务上提供高性能的共享存储,用户无 需修改应用程序,即可无缝迁移业务系统上云。

更多关于文件存储NAS的介绍,请参见文件存储NAS产品详情页。

● 容器服务 ACK

容器服务Kubernetes版(ACK)提供高性能可伸缩的容器应用管理能力,支持企业级容器化应用的全生命 周期管理。整合阿里云虚拟化、存储、网络和安全能力,打造云端最佳容器化应用运行环境。 更多关于容器服务ACK的介绍,请参见<mark>容器服务 ACK产品详情页</mark>。

# 详细信息

点击查看最佳实践详情

更多最佳实践

点击查看更多阿里云最佳实践