

K8s 发布说明

1. 概述

本文档包含了此次发布的 MXMACA K8s 软件包的特性，已知问题和使用限制等。

此次发布是 MXC500 GPU 的 **MXMACA-C500-K8s-0.8.1** 版本交付，适用于曦云®C500、C500X、C280、C290 OAM1.5、C290 OAM2.0、C550 OAM1.5、C550 OAM2.0、ARM 和曦思®N260。

表 1 列出了系统测试覆盖率和通过率。

表 1 系统测试覆盖率及通过率

系统	测试覆盖率/通过率
MXMACA-C500-K8s-0.8.1	Full regression with release quality

1.1 交付内容

此次发布的软件包包含以下内容：

- MXMACA-C500-K8s-0.8.1 K8s 软件栈

2. 新增特性及变更

本章列出历次发布的新增特性及变更。

2.1 MXMACA-C500-K8s-0.8.1

模块		特性说明
Extensions	volcano	volcano 提供了丰富的调度策略，增强型的 Job 管理能力及良好的生态支持
	gpu-aware	自定义权重，通过计算分数来进行资源调度

2.2 MXMACA-C500-K8s-0.8.0

模块		特性说明
Operator	driverpolicy	资源部署策略，支持选择使用 node driver 还是 container driver。
	用户容器镜像优化	MXMACA®由 GPU Operator 统一管理并部署在每个工作节点上，用户的作业运行时将使用节点上已部署的 MXMACA SDK。 此过程在操作层面对用户透明，尽管用户能够感知到容器在运行时自动安装了 MXMACA SDK 这一变化，但无需为此做任何额外操作。
	驱动自动部署	GPU Operator 方案提供了自动部署驱动及配置 GPU 虚拟化规格的能力。管理员可配置 <code>driver.deployPolicy</code> 选取不同部署策略，或关闭内核态驱动的自动部署功能。
	MXMACA 清理策略	Never: 不执行任何操作 OlderVersionFirst: 根据版本号排序，优先清理老版本 OlderTimestampFirst: 根据使用时间排序，优先清理未被使用的版本
Extensions	gpu-device	gpu-device 包含了资源分配器的逻辑实现，对于确定份数的 GPU 资源请求，gpu-device 总是确保从避免资源碎片，卡间互联拓扑角度进行最优分配。 gpu-device 会定期检查沐曦 GPU 设备的健康状态，识别出故障的 GPU 资源，并将其从可分配资源中移除。
	gpu-label	负责监控 k8s 节点上沐曦 GPU 及 MXMACA 软件栈的状态信息，并以标签的形式对节点进行标记。用户提交任务时，可通过在 <code>nodeSelector</code> 字段设置节点标签的形式来筛选符合预期的节点。

3. 已知问题和使用限制

模块	问题和限制说明
	暂无

飞腾信息 Metax Confidential
2024-11-18 15:00:00

声明

版权所有 ©2024 沐曦集成电路（上海）有限公司。保留所有权利。

本文档中呈现的信息属于沐曦集成电路（上海）有限公司和/或其附属公司（以下统称为“沐曦”），非经沐曦事先书面许可，任何实体或个人均不得获得本文档的副本，且无权以任何方式处理本文档，包括但不限于使用、复制、修改、合并、出版、发行、销售或传播本文档的部分或全部。

本文档内容仅供参考，不提供任何形式的、明示或暗示的保证，包括但不限于对适销性、适用于任何目的和/或不侵权的保证。在任何情况下，沐曦均不对因本文档引起的、由本文档造成的、或与之相关的任何索赔、损害或其他责任负责。

沐曦保留自行决定随时更改、修改、添加或删除本文档的部分或全部的权利。沐曦保留最终解释权。

沐曦、MetaX 和其他沐曦图标是沐曦的商标。本文档中提及的所有其他商标和商品名称均为其各自所有者的财产。

飞腾信息 MetaX Confidential
2024-11-18 15:00:00