

# UNIS DataEngine AIOS 人工智能引擎

## 产品概述

随着 AI 科学计算的技术突破，人工智能已经广泛应用于各行各业，如医疗、金融、汽车、法律、工业、教育等等，其中 AI 科学计算服务中心化也成为众望所归的选择，虽然 AI 深度学习目前的门槛有所降低，但是对于专业的数据科学家，依然是需要为了环境搭建，并行运算，分布式存储，作业调度等基础设施和服务耗费较多的人力和成本，为此，UNIS 公共科学计算 AIOS 平台应运而生，提出软硬件一体化方案，从基础硬件的部署和软件安装，到交互式开发环境的一键启动，从模型的深度训练和调优，到多机多卡 GPU 作业灵活调度，UNIS AIOS 平台，提供了十分简洁的使用方式，实现了资源的整合/弹性扩容缩容和合理调度，同时也提供丰富的可自定义的软件和镜像和二次开发的 API 接口，可方便的集成进入原有 SaaS 平台。

针对图像处理、语音识别、自然语言处理等深度学习场景下，需要搭建大规模的 GPU 集群，针对不同的算法模型、不同的深度学习框架，用户如何统一调度与管理 GPU 集群的计算资源、存储资源，分配给不同的租户使用，是首当其冲需要解决的问题。

对于 TensorFlow、Caffe、MxNet 等深度学习框架，如何快速部署，提供开发镜像环境，满足不同用户在不同场景下的框架需求、算法需求与开发需求，也是数据科学家难以逾越的一道门槛。

面对不同用户同时进行模型训练、在线推理，采用什么策略对各个任务进行调度，是抢占模式还是先进先出，以及每个训练任务利用哪个 GPU 加速卡，每个卡的运行状态如何，都需要统一的监控与管理。

针对以上问题，为用户提供一体化的软硬件部署和管理服务，减少开发者系统安装维护工作量；优化分布式训练部署模型，实现多机多卡 GPU 资源与训练作业灵活调度；提供丰富的可自定义软件和镜像库，充分满足客户对 AI 计算环境的需求。

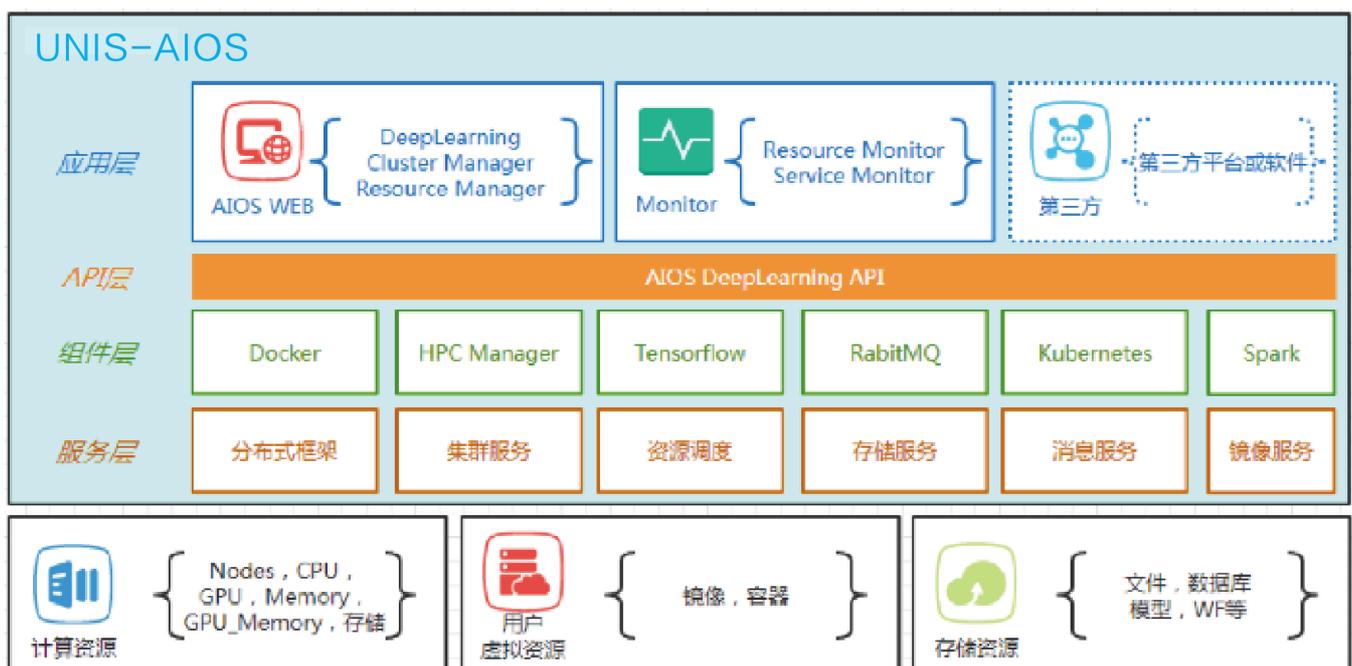
**灵活的资源调度机制：**提供强大的资源调度策略，以及资源实时监控，使企业可以有效、合理的使用各种计算资源。

**完善的 API 接口服务：**系统所有核心业务，都对外提供完善的 API 接口，用户可以通过这些接口，实现人工智能平台服务与用户已有 SaaS 平台的深度集成。

**灵活的权限管理策略：**系统通过对人员角色权限的划分，以及资源的使用规划，使得平台中不同的用户的计算资源都能很好的隔离，满足企业对权限管理的各种要求。

**丰富的性能监控服务：**AIOS 平台提供完善的性能监控服务，能实时监控系统所有服务的健康状况和硬件/网络利用率情况，并针对各种日常运维工作，提供可视化的操作界面，提高系统运维管理者的工作效率。

## 系统架构



### ◆ 先进的前后台分离架构

在以前传统的网站开发中，前端一般扮演的只是切图的工作，只是简单地将 UI 设计师提供的原型图实现成静态的 HTML 页面，而具体的页面交互逻辑，比如与后台的数据交互工作等，可能都是由后台的开发人员来实现的，或者是前端是紧紧的耦合后台。导致后台的开发压力大大增加，前后端工作分配不均。不仅仅开发效率慢，而且代码难以维护。

AIOS 采用先进的前后端分离架构，很好的解决前后端分工不均，开发过程相互依赖，bug 难以定位等诸多问题。将更多的用户交互逻辑由前端专职处理，而后端则可以专注于数据处理，业务权限控制等，前后端通过标准的 restful 接口实现数据交互。

后端专注于：服务层 & 数据访问层 & 权限控制；

前端专注于：页面展现（视图层）& 交互逻辑；

### ◆ 领先的微服务架构

AIOS 采用 kubernetes+docker+rabbitmq 的微服务架构模式，利用 kubernetes 实现高可用的集群环境，以及统一的资源调度，配合 docker 容器技术实现多租户资源隔离，由 rabbitmq 实现分布式消息处理。为平台展现层提供强大的内核支撑，平台采用的微服务架构模式有以下特点：

- ◎服务高度自治，集中管理；
- ◎复杂业务得到拆分，易于维护；
- ◎高度灵活易于拓展；

### ◆ 多租户存储资源隔离

平台使用 NFS (Network File System) 实现统一的网络文件存储系统，极大简化了平台部署的复杂性，提高了公共文件资源的利用率。再结合 linux 多用户多任务的系统特性，实现文件存储资源的多租户资源隔离。

### ◆ 强大的作业调度机制

AIOS 平台的核心是模型训练，对 CPU、内存、GPU 等资源的合理使用要求非常高，通常这些计算设备造价不菲，计算资源的最大化合理使用是体现一个计算平台最重要的指标。平台采用的 kubernetes 先天支持多种策略的作业调度，可以有效保证各类训练任务的及时有效执行。再结合平台提供的用户资源申请分配机制，以及资源使用率告警机制，可以灵活有效的管理多租户场景下，集群计算资源的统一合理调度。

### ◆ 完善的资源监控体系

Heapster 是容器集群监控和性能分析工具，可以定时采集集群环境中所有节点的 CPU、内存、网络以及磁盘情况，平台通过 Heapster 实现上述资源的统一采集及管理；而平台自主研发的 cMonitor 则可以对 GPU 资源进行定时的采集及管理；结合平台提供的计算资源告警机制，最终可以在平台展现层实时监控集群环境下的所有资源，并及时向管理员发送告警信息，有效管理集群硬件资源，提高资源利用率。。

## ➤ 产品主要特性及优势

AIOS 平台为用户提供了强大的全流程可视化管理平台：交互学习平台，集群管理平台，性能监控平台，审计平台这四大模块又同时对外提供丰富的接口组件，使得用户既可以完全通过我们的平台进行常规业务操作，又可以利用我们提供的接口组件，将核心服务集成到自己现有系统中。

AIOS 平台提供了一系列的函数库，方便用户在建模编写脚本的过程中，快速调用这些函数而无需关心这些业务无关的技术细节；同时我们在系统不同功能模块中，提供了不同的工作集，方便用户在具体场景中快速操作。所有这些，都大大加速了用户在开发人工智能解决方案时的速度。

### ◆ 统一的集群管理

负责整个系统计算资源的集中管理、统一分配与作业调度，包括 GPU 资源池的集中管理与分配、多租户方式隔离计算资源、以作业方式动态分配计算资源以及计算资源回收等。

### ◆ 统一的监控运维

实时监控管理集群资源使用情况和集群状态，包括作业状态、GPU 使用率、集群健康度等，并分析每一类的资源占用情况，提供触发预警机制。

### ◆ 统一的开发环境

提供一站式的交互开发操作界面，帮助用户完成模型脚本在线编辑、模型训练、模型验证以及模型推理等核心功能，并结合硬件资源可视化、作业调度器，最大化提高系统硬件资源的利用率。

## ➤ 产品功能特性

### ◆ 资源管理服务

系统提供完备的资源管理平台，对所有计算资源进行集中管理，通过该平台，可以实现对各类资源的状态查询以及相关维护操作。

源ID	源名称	IP	状态	节点	节点名称	CPU%	GPU%	GPU显存%	内存%	GPU利用率
1	cn-aios-01	10.2.1.10.134	空闲	运行	cn-aios-01					
2	cn-aios-02	10.2.1.10.135	空闲	运行	cn-aios-02					
3	cn-aios-03	10.2.1.10.136	空闲	运行	cn-aios-03					
4	cn-aios-04	10.2.1.10.137	空闲	运行	cn-aios-04					
5	cn-aios-05	10.2.1.10.138	空闲	运行	cn-aios-05					

### ◆ 作业调度服务

系统提供强大的作业调度引擎，为用户提供多种作业调度策略：先进先出，资源回填，公平共享，作业抢占，用户循环调度，用户作业均衡等。



## ◆ 性能监控服务

系统针对各种计算资源，提供多维度的资源性能监控指标查询接口，并图形化展示。



## ◆ API 服务

◎深度学习模块核心服务如下：

- ◆ 文件服务
- ◆ 数据集服务
- ◆ 模型存储服务
- ◆ 模型实例化服务
- ◆ 训练服务
- ◆ 超参搜索服务
- ◆ 评估服务
- ◆ 推理服务
- ◆ 指标监控服务

另外针对 AI 深度学习训练过程中复杂不可解释的参数选配，为了减少尝试次数，浪费资源和时间，也提供丰富的超参算法支持，支持 Random Search、TPE(Tree-based Parzen Estimator)以及 Bayesian 超参搜索算法，利用高效的超参搜索算法实现并行超参搜索，充分发挥集群计算能力，多任务并发搜索，不同任务间分享搜索结果（不同的），这样以改进效率为目标，做到搜索效果与搜索代价的良好平衡，还有提供蒙特卡洛树搜索 + 深度学习网络功能，解决搜索空间过大问题，并对搜索结果进行学习。

## 运行环境

### ◆ 计算节点所需硬件配置

指标项	最低配置	推荐配置
机器数量	1	2(支持扩展)
型号	X86 平台的服务器	X86 平台的服务器
CPU	32 核（物理核数，非超线程核数），支持 AVX 模式	32 核（物理核数，非超线程核数）或更高，支持 AVX 模式
GPU	n*Tesla P4	n*Tesla P40/P100/V100
内存	128G 内存	256G 以上内存
磁盘	2*1T 硬盘作 RAID1 磁盘阵列	2*2T 以上容量硬盘作 RAID1 磁盘阵列
网卡	1 个千兆网卡	1 个万兆以上网卡

### ◆ 选配信息

项目	描述
大数据平台部署服务(4 台)	必配
UNIS 公共科学计算 AIOS 软件 License 费用	必配, 3 个节点
UNIS 大数据技术支持服务(一年)	必配



北京紫光恒越网络科技有限公司

北京基地  
北京市海淀区中关村东路 1 号院 2 号楼 402 室  
邮编: 100084  
电话: 010-62166890  
传真: 010-51652020-116  
版本:

Copyright ©2012 北京紫光恒越网络科技有限公司 保留一切权利

免责声明: 虽然 UNIS 试图在本资料中提供准确的信息, 但不保证资料的内容不含有技术性误差或印刷性错误, 为此 UNIS 对本资料中的不准确不承担任何责任。  
UNIS 保留在没有通知或提示的情况下对本资料的内容进行修改的权利。

<http://www.unishy.com>

客户服务热线  
**400-910-9998**