

盘古大模型
3.3.0

产品介绍

文档版本 01
发布日期 2025-02-27



版权所有 © 华为云计算技术有限公司 2025。保留一切权利。

非经本公司书面许可，任何单位和个人不得擅自摘抄、复制本文档内容的部分或全部，并不得以任何形式传播。

商标声明



HUAWEI和其他华为商标均为华为技术有限公司的商标。

本文档提及的其他所有商标或注册商标，由各自的所有人拥有。

注意

您购买的产品、服务或特性等应受华为云计算技术有限公司商业合同和条款的约束，本文档中描述的全部或部分产品、服务或特性可能不在您的购买或使用范围之内。除非合同另有约定，华为云计算技术有限公司对本文档内容不做任何明示或暗示的声明或保证。

由于产品版本升级或其他原因，本文档内容会不定期进行更新。除非另有约定，本文档仅作为使用指导，本文档中的所有陈述、信息和建议不构成任何明示或暗示的担保。

华为云计算技术有限公司

地址：贵州省贵安新区黔中大道交兴功路华为云数据中心 邮编：550029

网址：<https://www.huaweicloud.com/>

目录

1 什么是盘古大模型	1
2 产品优势	2
3 应用场景	3
4 产品功能	4
4.1 空间管理	4
4.2 数据工程	5
4.3 模型开发	6
4.4 Agent 开发	6
5 模型能力与规格	8
5.1 盘古 NLP 大模型能力与规格	8
5.2 盘古科学计算大模型能力与规格	12
5.3 盘古专业大模型能力与规格	16
6 基础知识	18
6.1 大模型开发基本流程介绍	18
6.2 大模型开发基本概念	19
7 安全	22
7.1 责任共担	22
7.2 身份认证与访问控制	23
7.3 数据保护技术	23
7.4 审计	23
8 权限管理	25
9 约束与限制	28
10 与其他服务的关系	30

1 什么是盘古大模型

盘古大模型服务致力于深耕行业，打造多领域行业大模型和能力集。ModelArts Studio大模型开发平台是盘古大模型服务推出的集数据管理、模型训练和模型部署为一体的一站式大模型开发平台及大模型应用开发平台，盘古NLP大模型、科学计算大模型、专业大模型能力通过ModelArts Studio大模型开发平台承载，它提供了包括盘古大模型在内的多种大模型服务，提供覆盖全生命周期的大模型工具链。

ModelArts Studio大模型开发平台为开发者提供了一种简单、高效的开发和部署大模型的方式。平台提供了包括数据处理、模型训练、模型部署、Agent开发等功能，以帮助开发者充分利用盘古大模型的功能。企业可以根据自己的需求选取合适的大模型相关服务和产品，方便地构建自己的模型和应用。

- **数据工程工具链**

数据是大模型训练的基础，为大模型提供了必要的知识和信息。数据工程工具链作为盘古大模型服务的重要组成部分，具备数据获取、清洗、数据合成、数据标注、数据评估、数据配比、数据发布和管理等功能。

该工具链能够高效收集和處理各种格式的数据，满足不同训练和评测任务的需求。通过提供自动化的质量检测和数据清洗能力，对原始数据进行优化，确保其质量和一致性。同时，数据工程工具链还提供强大的数据存储和管理能力，为大模型训练提供高质量的数据支撑。

- **模型开发工具链**

模型开发工具链是盘古大模型服务的核心组件，提供从模型创建到部署的一站式解决方案。

该工具链具备模型训练、压缩、部署、评测、推理等功能，通过高效的推理性能和跨平台迁移工具，模型开发工具链能够保障模型在不同环境中的高效应用。

- **Agent开发工具链**

应用开发工具链是盘古大模型平台的关键模块，支持提示词工程和智能Agent应用创建。该工具链提供提示词设计和管理工具，优化大模型的输入提示，提升输出的准确性和相关性。通过可视化编排工具，应用开发工具链加速大模型应用的开发，满足复杂业务需求。

2 产品优势

- **预置多，数据工程“易”**
ModelArts Studio大模型开发平台预置多种数据处理AI算子，多种标注工具，满足用户多任务多场景需求，提高开发/标注效率>10X。
- **0代码，模型开发“简”**
ModelArts Studio大模型开发平台预置盘古系列预训练大模型，支持快速开发，全程0代码开发，极大降低大模型开发门槛。
- **功能强，Agent开发“好”**
Agent开发提供便捷搭建大模型应用功能，并提供功能强大的插件配置，让Agent能力更强，更专业。
- **统一管，资产管理“全”**
ModelArts Studio大模型开发平台数据、模型、Agent应用在统一的入口进行管理，可以快速的掌握资产的使用情况、版本情况和溯源信息等。

3 应用场景

客服

通过NLP大模型对传统的客服系统进行智能化升级，提升智能客服的效果。企业原智能客服系统仅支持回复基础的FAQ，无语义泛化能力，意图理解能力弱，转人工频率极高。面对活动等时效性场景，智能客服无回答能力。提高服务效率：大模型智能客服可以7x24小时不间断服务，相较于人工客服，可以处理更多的客户咨询，且响应速度快；降低运营成本：企业可以通过智能客服处理大部分的常规问题，将人工客服释放出来处理更复杂、更个性化的客户需求；个性化服务：基于大模型的智能客服能够学习和适应用户的行为模式和偏好，提供更加个性化的服务。

农业

科学计算大模型包括全球中期天气要素模型和降水模型，可以对未来一段时间的天气和降水进行预测，全球中期天气要素模型和降水模型能够在全全球范围内进行预测，不仅仅局限于某个地区。它的分辨率相当于赤道附近每个点约25公里x25公里的空间。通过降水模型预测未来的降雨情况，农民和农业管理者可以更有效地规划灌溉时间和频率，也能为可能发生的干旱提供预警，使农业部门能够及时采取措施，如推广节水技术或调整种植计划。

代码助手

在软件开发领域，编程语言的多样性和复杂性给程序员带来了巨大的挑战。盘古NLP大模型为程序员提供了强大的代码助手，显著提升了研发效率。

盘古大模型能够根据用户给定的题目，快速生成高质量的代码，支持Java、Python、Go等多种编程语言。它不仅能够提供完整的代码实现，还能够根据用户的需求，进行代码补全和不同编程语言之间的改写转化。借助盘古大模型，程序员可以更加专注于创新和设计，而无需过多关注繁琐的编码工作。它不仅提升了代码的质量和稳定性，还缩短了开发周期，加速了产品的迭代和发布。

4 产品功能

4.1 空间管理

ModelArts Studio大模型开发平台为用户提供了灵活且高效的空间资产管理方式。平台支持用户根据不同的使用场景、项目类别或团队需求，自定义创建多个工作空间。每个工作空间都是完全独立的，确保了工作空间内的资产不受其他空间的影响，从而保障数据和资源的隔离性与安全性。用户可以根据需求灵活划分工作空间，实现资源的有序管理与优化配置，确保各类资源在不同场景中的最大化利用。为进一步优化资源的管理，平台还提供了多种角色权限体系。用户可以根据自身角色从管理者到各模块人员进行不同层级的权限配置，确保每个用户在其指定的工作空间内，拥有合适的访问与操作权限。这种精细化的权限管理方式，既保证了数据的安全性，又提高了资源的高效利用。

在平台中，空间资产指的是存储在工作空间中的所有资源，包括数据资产和模型资产。这些资产是用户在平台上进行开发和管理的基礎，集中存储和统一管理的方式有助于提升操作效率，并确保资源的规范性与安全性。

- **数据资产：**数据资产是指用户在平台上发布的所有数据集。这些数据集会被存储在数据资产中，用户可以随时查看数据集的详细信息，如数据格式、大小、配比比例等，同时平台会自动记录每个数据集的操作历史，例如创建、发布及上线过程。为了进一步简化管理，平台还支持数据集的删除功能，使用户能够对数据集进行灵活管理和调整。在模型训练和数据分析过程中，用户可以根据需求调用这些数据集，确保数据的准确性与安全性，从而提升数据资产的利用率。同时支持数据集发布到Gallery，支持从Gallery订阅数据集。
- **模型资产：**模型资产包括用户试用、订购或在平台上训练后发布的模型，这些模型统一存储在模型资产中，便于集中管理。用户可以查看模型的所有历史版本及操作记录，从而了解模型的演变过程。同时，平台支持一系列便捷的模型操作，如模型训练、压缩和部署，帮助用户简化模型开发和应用流程。此外，平台还提供了导入和导出功能，支持用户将其他局点的盘古大模型迁移到本地局点，这使得模型资产在不同局点间的共享和管理变得更加灵活高效。同时支持模型发布到Gallery，支持从Gallery订阅模型。

通过统一管理空间资产，平台不仅帮助用户高效组织和利用资源，还保障了资产的安全性、一致性与灵活性。这些功能的结合，确保了平台上资源的高效利用与智能配置，为用户提供了更为便捷的开发和管理体验。

4.2 数据工程

ModelArts Studio开发平台提供了全面的数据工程功能。该模块涵盖数据获取、加工、标注、评估和发布等关键环节，帮助用户高效构建高质量的训练数据集，推动AI应用的成功落地。具体功能如下：

- **数据获取：**用户可以轻松将多种类型的数据导入ModelArts Studio大模型开发平台，支持的数据类型包括文本、图片、视频、气象、预测数据以及用户自定义的其他类型数据。平台提供灵活的数据接入方式以及支持多种文件格式导入，确保不同业务场景下的数据获取需求得到满足。
- **数据加工：**平台提供强大的数据加工功能，可以对文本、视频、图片、气象类型的数据进行数据提取、过滤、转换、打标签和评分等加工处理。针对不同类型的数据集，平台提供了专用的清洗算子以及支持用户创建自定义算子实现个性化的数据清洗诉求。确保生成高质量的训练数据以满足业务需求和模型训练的要求。用户还可以灵活地调整算子编排顺序以及自定义清洗模板，有效提升数据清洗效率并支持大规模数据处理，确保生成的数据集符合训练的标准。
- **数据合成：**平台支持利用预置或自定义的数据指令对预训练文本、单轮问答、单轮问答（人设）数据集类型进行处理，并根据设定的轮数生成新数据。通过数据合成技术，可以生成大量高质量的训练数据，这些数据可以用于大模型的预训练，增强模型的泛化能力和性能。
- **数据标注：**平台支持对无标签的数据添加标签或对现有的标签进行重新标注，以提升数据集的标注质量。用户可以针对不同的数据集灵活地选择对应的标注项，还可以自定义选择多人标注、审核以及标注任务移交。针对文本和图片类数据集，平台还提供AI预标注功能。利用盘古大模型的智能能力，显著降低人工标注的工作量和成本，从而显著地提高标注效率。
- **数据评估：**平台支持对处理后的文本、图片、视频等多种格式数据进行质量评估，并预置了基础的评估标准，用户可以直接使用预置标准或创建自定义评估标准，以满足个性化的数据质量需求。最终生成详细的质量评估报告，这些报告能够帮助用户检验数据的准确性、完整性和一致性，确保数据在进行模型训练前的高质量标准，以保证模型在实际应用中的可靠性和稳定性。
- **数据配比：**平台支持对文本、图片类数据进行数据配比。用户在勾选数据集时可以勾选多条，通过调整不同来源或类型数据的比例，以优化模型训练过程。通过数据配比可以确保模型能够更全面地学习和理解数据的多样性，提高模型的泛化能力和性能。
- **数据发布：**平台支持数据集发布。用户可以将处理后的数据集发布为多种格式，包括标准格式和盘古格式。尤其对于文本类和图片类数据集，平台支持将其转换为专门用于训练盘古大模型的盘古格式，为后续模型训练提供高效的数据支持。
- **数据管理：**平台支持数据全链路血缘追溯，用户单击数据集名称可以在“数据血缘”页签，查看该数据集所经历的操作。全链路血缘追溯可以帮助用户正向实现数据集影响分析，逆向实现快速问题追踪，提升数据运维和数据治理的效率，帮助用户更好地对数据进行追根溯源。另外平台还提供了完善的标签体系、支持数据按行业标准进行分类、按行业标准进行安全分级、内置场景分类标签。帮助用户进行数据分类、数据质量控制和数据资产管理，提升数据治理的效率和效果。

通过整合上述功能，数据工程在AI研发中不仅帮助用户高效构建高质量的训练数据集，还通过全流程的数据处理和管理，探索数据与模型性能的内在联系，为模型训练和应用提供坚实的数据基础，推动了模型的精确训练与持续优化，提升了AI应用开发的效率和成果的可靠性。

4.3 模型开发

ModelArts Studio大模型开发平台提供了模型开发功能，涵盖了从模型训练到模型调用的各个环节。平台支持全流程的模型生命周期管理，确保从数据准备到模型部署的每一个环节都能高效、精确地执行，为实际应用提供强大的智能支持。

- **模型训练**：在模型开发的第一步，ModelArts Studio大模型开发平台为用户提供了丰富的训练工具与灵活的配置选项。用户可以根据实际需求选择合适的模型架构，并结合不同的训练数据进行精细化训练。平台支持分布式训练，能够处理大规模数据集，从而帮助用户快速提升模型性能。该模块提供预训练、全量微调、LoRA微调等。
- **模型评测**：为了确保模型的实际应用效果，平台提供了多维度的模型评测功能。通过自动化的评测机制，用户可以在训练过程中持续监控模型的精度、召回率等关键指标，及时发现潜在问题并优化调整。评测功能能够帮助用户在多种应用场景下验证模型的准确性与可靠性。支持基于规则的自动评测方式，NLP模型展示准确率，F1分数，BLEU、ROUGE等自动评测指标，支持支持人工评测自定义配置评测指标；并且支持基于人工评价操作界面，对模型表现从不同评价指标进行打分。
- **模型压缩**：在模型部署前，进行模型压缩是提升推理性能的关键步骤。通过压缩模型，能够有效减少推理过程中的显存占用，节省推理资源，同时提高计算速度。当前，平台支持对NLP大模型进行压缩，目前支持INT8、INT4量化压缩。
- **模型部署**：平台提供了一键式模型部署功能，用户可以轻松将训练好的模型部署到云端或本地环境中。平台支持多种部署模式，能够满足不同场景的需求。通过灵活的API接口，模型可以无缝集成到各类应用中。
- **模型调用**：在模型部署后，用户可以通过模型调用功能快速访问模型的服务。平台提供了高效的API接口，确保用户能够方便地将模型嵌入到自己的应用中，实现智能对话、文本生成等功能。

4.4 Agent 开发

Agent开发平台为开发者提供了一个全面的工具集，帮助您高效地开发、优化和部署应用智能体。无论您是新手还是有经验的开发者，都能通过平台提供的提示词工程、插件扩展、灵活的工作流设计和全链路调测功能，快速实现智能体应用的开发与落地，加速行业AI应用的创新与应用。

- **对于零码开发者（无代码开发经验的用户）：**
 - 平台提供了Prompt提示词工程和插件自定义等功能，帮助用户在无需编写代码的情况下，快速构建、调优并运行属于自己的大模型应用。通过简单的配置，用户可以轻松创建Agent应用，快速体验智能化应用的便捷性。
 - 平台提供导入知识功能，支持用户存储和管理数据，并与AI应用进行互动。支持多种格式的本地文档（如docx、pptx、pdf等），方便导入至知识，为Agent应用提供个性化数据支持。
 - 平台还提供全链路信息观测和调试工具，支持开发者深入分析Agent执行过程中的每个环节。通过对信息进行分层展示，帮助开发者优化AI应用的性能和稳定性，确保应用在不同环境下的顺畅运行。
- **对于低码开发者（具有一定代码开发经验的用户）：**
 - 基于上述功能，平台还提供了灵活的工作流设计功能，支持用户编写少量代码来构建逻辑复杂、稳定性要求高的Agent应用。通过拖拉拽方式，开发者可

- 以组合各种组件（如大模型、代码、意图识别等），快速搭建 workflows，实现更高效的应用开发。
- 平台还提供全链路信息观测和调试工具，支持开发者深入分析 workflow 执行过程中的每个环节。通过对信息进行分层展示，帮助开发者优化 AI 应用的性能和稳定性，确保应用在不同环境下的顺畅运行。

5 模型能力与规格

5.1 盘古 NLP 大模型能力与规格

盘古NLP大模型是业界首个超千亿参数的中文预训练大模型，结合了大数据预训练和多源知识，借助持续学习不断吸收海量文本数据，持续提升模型性能。除了实现行业知识检索、文案生成、阅读理解等基础功能外，盘古NLP大模型还具备模型调用等高级特性，可在智能客服、创意营销等多个典型场景中，提供强大的AI技术支持。

ModelArts Studio大模型开发平台为用户提供了多种规格的NLP大模型，以满足不同场景和需求。不同模型在处理上下文token长度和功能上有所差异，以下是当前支持的模型清单，您可以根据实际需求选择最合适的模型进行开发和应用。

表 5-1 盘古 NLP 大模型规格

模型支持区域	模型名称	可处理最大上下文长度	可处理最大输出长度	说明
西南-贵阳	Pangu-NLP-N1-Chat-32K-20241130	32K	4K	2024年11月发布的版本，支持8K序列长度训练，4K/32K序列长度推理。全量微调、LoRA微调8个训练单元起训，1个推理单元即可部署。
	Pangu-NLP-N1-Chat-128K-20241130	128K	4K	2024年11月发布的版本，仅支持128K序列长度推理。
	Pangu-NLP-N1-32K-3.1.34	32K	4K	2024年11月发布的版本，支持8K序列长度训练，4K/32K序列长度推理。全量微调、LoRA微调8个训练单元起训，1个推理单元即可部署，4K支持256并发，32K支持256并发。

模型支持区域	模型名称	可处理最大上下文长度	可处理最大输出长度	说明
	Pangu-NLP-N1-32K-3.2.36	32K	4K	2025年1月发布的版本，支持32K序列长度训练，4K/32K序列长度推理。全量微调、LoRA微调8个训练单元起训，1个推理单元即可部署，4K支持256并发，32K支持256并发。
	Pangu-NLP-N1-128K-3.1.34	128K	4K	2024年11月发布的版本，仅支持128K序列长度推理，4卡2并发。
	Pangu-NLP-N1-128K-3.2.36	128K	4K	2025年1月发布的版本，仅支持128K序列长度推理，4个推理单元8并发。
	Pangu-NLP-N2-Base-20241030	-	4K	2024年11月发布的版本，仅支持模型增量预训练。32个训练单元起训，预训练后的模型版本需要通过微调之后，才可支持推理部署。
	Pangu-NLP-N2-Chat-32K-20241030	32K	4K	2024年10月发布版本，支持8K序列长度训练，4K/32K序列长度推理。全量微调32个训练单元起训，LoRA微调8个训练单元起训，4个推理单元即可部署。此模型版本差异化支持预训练特性、INT8量化特性。
	Pangu-NLP-N2-4K-3.2.35	4K	4K	2025年1月发布的版本，支持4K序列长度训练，4K序列长度推理。全量微调32个训练单元起训，LoRA微调8个训练单元起训，4个推理单元即可部署，支持192并发。此模型版本差异化支持预训练特性、INT8量化特性。
	Pangu-NLP-N2-32K-3.1.35	32K	4K	2024年12月发布版本，支持8K序列长度训练，4K/32K序列长度推理。全量微调32个训练单元起训，LoRA微调8个训练单元起训，4个推理单元即可部署，4K支持64并发，32K支持64并发。此模型版本差异化支持预训练特性、INT8量化特性。
	Pangu-NLP-N2-32K-3.1.35	32K	4K	2025年1月发布的版本，支持32K序列长度训练，32K序列长度推理。全量微调32个训练单元起训，LoRA微调8个训练单元起训，4个推理单元即可部署，支持128并发。此模型版本差异化支持预训练特性、INT8量化特性。

模型支持区域	模型名称	可处理最大上下文长度	可处理最大输出长度	说明
	Pangu-NLP-N2-128K-3.1.35	128K	4K	2024年12月发布的版本，仅支持128K序列长度推理部署，8个推理单元64并发。
	Pangu-NLP-N2-256K-3.1.35	256K	4K	2024年12月发布的版本，仅支持256K序列长度推理部署，8个推理单元64并发。
	Pangu-NLP-N4-Chat-4K-20241130	32K	4K	2024年11月发布的版本，支持4K序列长度训练，4K序列长度推理。全量微调64个训练单元起训，LoRA微调32个训练单元起训，8个训练单元即可部署。此模型版本差异化支持预训练、INT8/INT4量化特性。
	Pangu-NLP-N4-Chat-32K-20241130	32K	4K	2024年11月发布的版本，仅支持32K序列长度推理部署。
	Pangu-NLP-N4-4K-2.5.32	4K	4K	2024年11月发布的版本，支持4K序列长度训练，4K序列长度推理。全量微调64个训练单元起训，LoRA微调32个训练单元起训，8个推理单元即可部署，支持64并发。此模型版本差异化支持预训练、INT8/INT4量化特性。
	Pangu-NLP-N4-4K-2.5.35	4K	4K	2025年1月发布的版本，支持4K序列长度训练，4K序列长度推理。全量微调64个训练单元起训，LoRA微调32个训练单元起训，8个推理单元即可部署，支持128并发。此模型版本差异化支持预训练、INT8/INT4量化特性。
	Pangu-NLP-N4-32K-2.5.32	32K	4K	2024年11月发布的版本，仅支持32K序列长度推理部署，8个推理单元64并发。
	Pangu-NLP-N4-32K-2.5.35	32K	4K	2025年1月发布的版本，仅支持32K序列长度推理部署，8个推理单元128并发。

在选择和使用盘古大模型时，了解不同模型所支持的操作行为至关重要。不同模型在预训练、微调、模型评测、在线推理和能力调测等方面的支持程度各不相同，开发者应根据自身需求选择合适的模型。以下是盘古NLP大模型支持的具体操作：

表 5-2 盘古 NLP 大模型支持的能力

模型	预训练	微调	模型评测	模型压缩	在线推理	能力调测
Pangu-NLP-N1-Chat-32K-20241130	-	√	√	-	√	√
Pangu-NLP-N1-Chat-128K-20241130	-	-	√	-	√	√
Pangu-NLP-N1-32K-3.1.34	-	√	√	-	√	√
Pangu-NLP-N1-32K-3.2.36	-	√	√	√	√	√
Pangu-NLP-N1-128K-3.1.34	-	-	√	-	√	√
Pangu-NLP-N1-128K-3.2.36	-	-	√	-	√	√
Pangu-NLP-N2-Base-20241030	√	-	-	-	-	-
Pangu-NLP-N2-Chat-32K-20241030	-	√	√	√	√	√
Pangu-NLP-N2-4K-3.2.35	√	√	√	√	√	√
Pangu-NLP-N2-32K-3.1.35	√	√	√	√	√	√
Pangu-NLP-N2-128K-3.1.35	-	-	√	-	√	√
Pangu-NLP-N2-256K-3.1.35	-	-	√	-	√	√
Pangu-NLP-N4-Chat-4K-20241130	-	√	√	√	√	√
Pangu-NLP-N4-Chat-32K-20241130	-	√	√	√	√	√
Pangu-NLP-N4-4K-2.5.32	√	√	√	√	√	√
Pangu-NLP-N4-4K-2.5.35	√	√	√	√	√	√

模型	预训练	微调	模型评测	模型压缩	在线推理	能力调测
Pangu-NLP-N4-32K-2.5.32	-	-	√	-	√	√
Pangu-NLP-N4-32K-2.5.35	-	-	√	-	√	√

5.2 盘古科学计算大模型能力与规格

盘古科学计算大模型面向气象、医药、水务、机械、航天航空等领域，融合了AI数据建模和AI方程求解方法。该模型从海量数据中提取数理规律，利用神经网络编码微分方程，通过AI模型更快速、更精准地解决科学计算问题。

ModelArts Studio大模型开发平台为用户提供了多种规格的科学计算大模型，以满足不同场景和需求。以下是当前支持的模型清单，您可以根据实际需求选择最合适的模型进行开发和应用。

表 5-3 盘古科学计算大模型规格

模型支持区域	模型名称	说明
西南-贵阳一	Pangu-AI4S-Ocean_24h-20241130	2024年11月发布的版本，用于海洋基础要素预测，可支持1个实例部署推理。
	Pangu-AI4S-Ocean_24h-3.1.0	2025年1月发布的版本，用于海洋基础要素预测，可支持1个实例部署推理。
	Pangu-AI4S-Ocean_Regional_24h-20241130	2024年11月发布的版本，用于区域海洋基础要素预测，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Ocean-Regional_24h-3.1.0	2025年1月发布的版本，用于区域海洋基础要素预测，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Ocean_Ecology_24h-20241130	2024年11月发布的版本，用于海洋生态要素预测，可支持1个实例部署推理。
	Pangu-AI4S-Ocean-Ecology_24h-3.1.0	2025年1月发布的版本，用于海洋生态要素预测，可支持1个实例部署推理。
	Pangu-AI4S-Ocean_Swell_24h-20241130	2024年11月发布的版本，用于海浪预测，可支持1个实例部署推理。
	Pangu-AI4S-Ocean-Swell_24h-3.1.0	2025年1月发布的版本，用于海浪预测，可支持1个实例部署推理。

模型支持区域	模型名称	说明
	Pangu-AI4S-Weather_Precip-20241030	2024年10月发布的版本，用于降水预测，支持1个实例部署推理。
	Pangu-AI4S-Weather-Precip_6h-3.0.0	2024年12月发布的版本，相较于10月发布的版本模型运行速度有提升，用于降水预测，支持1个实例部署推理。
	Pangu-AI4S-Weather-Precip_6h-3.1.0	2025年1月发布的版本，用于降水预测，支持1个实例部署推理。
	Pangu-AI4S-Weather_1h-20241030	2024年10月发布的版本，用于天气基础要素预测，时间分辨率为1小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_1h-3.0.0	2024年12月发布的版本，相较于10月发布的版本模型运行速度有提升，用于天气基础要素预测，时间分辨率为1小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_1h-3.1.0	2025年1月发布的版本，用于天气基础要素预测，时间分辨率为1小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_3h-20241030	2024年10月发布的版本，用于天气基础要素预测，时间分辨率为3小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_3h-3.0.0	2024年12月发布的版本，相较于10月发布的版本模型运行速度有提升，用于天气基础要素预测，时间分辨率为3小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_3h-3.1.0	2025年1月发布的版本，用于天气基础要素预测，时间分辨率为3小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_6h-20241030	2024年10月发布的版本，用于天气基础要素预测，时间分辨率为6小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_6h-3.0.0	2024年12月发布的版本，用于天气基础要素预测，时间分辨率为6小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_6h-3.1.0	2025年1月发布的版本，用于天气基础要素预测，时间分辨率为6小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_6h-3.1.1	2025年1月发布的版本，用于天气基础要素预测，时间分辨率为6小时，相较于3.1.0版本预报准确度更高，1个实例部署。

模型支持区域	模型名称	说明
	Pangu-AI4S-Weather_24h-20241030	2024年10月发布的版本，用于天气基础要素预测，时间分辨率为24小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_24h-3.0.0	2024年12月发布的版本，相较于10月发布的版本运行速度有提升，用于天气基础要素预测，时间分辨率为24小时，1个训练单元起训及1个实例部署。
	Pangu-AI4S-Weather_24h-3.1.0	2025年1月发布的版本，用于天气基础要素预测，时间分辨率为24小时，1个训练单元起训及1个实例部署。

在选择和使用盘古大模型时，了解不同模型所支持的操作行为至关重要。不同模型在预训练、微调、模型评测、模型压缩、在线推理和能力调测等方面的支持程度各不相同，开发者应根据自身需求选择合适的模型。以下是盘古科学计算大模型支持的具体操作：

表 5-4 盘古科学计算大模型支持的操作

模型	预训练	微调	模型评测	模型压缩	在线推理	能力调测
Pangu-AI4S-Ocean_24h-20241130	-	-	-	-	√	√
Pangu-AI4S-Ocean_24h-3.1.0	-	-	-	-	√	√
Pangu-AI4S-Ocean_Regional_24h-20241130	√	√	-	-	√	√
Pangu-AI4S-Ocean-Regional_24h-3.1.0	√	√	-	-	√	√
Pangu-AI4S-Ocean_Ecology_24h-20241130	-	-	-	-	√	√
Pangu-AI4S-Ocean-Ecology_24h-3.1.0	-	-	-	-	√	√

模型	预训练	微调	模型评测	模型压缩	在线推理	能力调测
Pangu-AI4S-Ocean_Swell_24h-20241130	-	-	-	-	√	√
Pangu-AI4S-Ocean-Swell_24h-3.1.0	-	-	-	-	√	√
Pangu-AI4S-Weather_Precip-20241030	-	-	-	-	√	√
Pangu-AI4S-Weather-Precip_6h-3.0.0	-	-	-	-	√	√
Pangu-AI4S-Weather-Precip_6h-3.1.0	-	-	-	-	√	√
Pangu-AI4S-Weather_1h-20241030	√	√	-	-	√	√
Pangu-AI4S-Weather_1h-3.0.0	√	√	-	-	√	√
Pangu-AI4S-Weather_1h-3.1.0	√	√	-	-	√	√
Pangu-AI4S-Weather_3h-20241030	√	√	-	-	√	√
Pangu-AI4S-Weather_3h-3.0.0	√	√	-	-	√	√
Pangu-AI4S-Weather_3h-3.1.0	√	√	-	-	√	√
Pangu-AI4S-Weather_6h-20241030	√	√	-	-	√	√
Pangu-AI4S-Weather_6h-3.0.0	√	√	-	-	√	√
Pangu-AI4S-Weather_6h-3.1.0	√	√	-	-	√	√
Pangu-AI4S-Weather_6h-3.1.1	-	-	-	-	√	√

模型	预训练	微调	模型评测	模型压缩	在线推理	能力调测
Pangu-AI4S-Weather_24h-20241030	√	√	-	-	√	√
Pangu-AI4S-Weather_24h-3.0.0	√	√	-	-	√	√
Pangu-AI4S-Weather_24h-3.1.0	√	√	-	-	√	√

5.3 盘古专业大模型能力与规格

盘古专业大模型是盘古百亿级NL2SQL模型，适用于问数场景下的自然语言问题到SQL语句生成，支持常见的聚合函数（如去重、计数、平均、最大、最小、合计）、分组、排序、比较、条件（逻辑操作、离散条件、范围区间等条件的混合和嵌套）、日期操作，支持多表关联查询。

与非专业大模型相比，专业大模型针对特定场景优化，更适合执行数据分析、报告生成和业务洞察等任务。

ModelArts Studio大模型开发平台为用户提供了多种规格的专业大模型，以满足不同场景和需求。以下是当前支持的模型清单，您可以根据实际需求选择最合适的模型进行开发和应用。

模型支持区域	模型名称	说明
西南-贵阳一	Pangu-NLP-BI-4K-20241130	2024年11月发布的版本，支持4K序列长度推理，支持4个推理单元部署。
	Pangu-NLP-BI-32K-20241130	2024年11月发布的版本，支持32K序列长度推理，支持8个推理单元部署。

在选择和使用盘古大模型时，了解不同模型所支持的操作行为至关重要。不同模型在预训练、微调、模型压缩、在线推理和能力调测等方面的支持程度各不相同，开发者应根据自身需求选择合适的模型。以下是盘古专业大模型支持的具体操作：

模型	预训练	微调	模型压缩	在线推理	能力调测
Pangu-NLP-BI-4K-20241130	-	-	-	√	-
Pangu-NLP-BI-32K-20241130	-	-	-	√	-

6 基础知识

6.1 大模型开发基本流程介绍

大模型（Large Models）通常指的是具有海量参数和复杂结构的深度学习模型，广泛应用于自然语言处理（NLP）等领域。开发一个大模型的流程可以分为以下几个主要步骤：

- **数据集准备**：大模型的性能往往依赖于大量的训练数据。因此，数据集准备是模型开发的第一步。首先，需要根据业务需求收集相关的原始数据，确保数据的覆盖面和多样性。例如，若是自然语言处理任务，可能需要大量的文本数据；如果是计算机视觉任务，则需要图像或视频数据。

- **数据预处理**：数据预处理是数据准备过程中的重要环节，旨在提高数据质量和适应模型的需求。常见的数据预处理操作包括：

- 去除重复数据：确保数据集中每条数据的唯一性。
- 填补缺失值：填充数据中的缺失部分，常用方法包括均值填充、中位数填充或删除缺失数据。
- 数据标准化：将数据转换为统一的格式或范围，特别是在处理数值型数据时（如归一化或标准化）。
- 去噪处理：去除无关或异常值，减少对模型训练的干扰。

数据预处理的目的是保证数据集的质量，使其能够有效地训练模型，并减少对模型性能的不利影响。

- **模型开发**：模型开发是大模型项目中的核心阶段，通常包括以下步骤：

- 选择合适的模型：根据任务目标选择适当的模型。
- 模型训练：使用处理后的数据集训练模型。
- 超参数调优：选择合适的学习率、批次大小等超参数，确保模型在训练过程中能够快速收敛并取得良好的性能。

开发阶段的关键是平衡模型的复杂度和计算资源，避免过拟合，同时保证模型能够在实际应用中提供准确的预测结果。

- **应用与部署**：当大模型训练完成并通过验证后，进入应用阶段。主要包括以下几个方面：

- 模型优化与部署：将训练好的大模型部署到生产环境中，可能通过云服务或本地服务器进行推理服务。此时要考虑到模型的响应时间和并发能力。

- 模型监控与迭代：部署后的模型需要持续监控其性能，并根据反馈进行定期更新或再训练。随着新数据的加入，模型可能需要进行调整，以保证其在实际应用中的表现稳定。

在应用阶段，除了将模型嵌入到具体业务流程中外，还需要根据业务需求不断对模型进行优化，使其更加精准和高效。

6.2 大模型开发基本概念

大模型相关概念

概念名	说明
大模型是什么	大模型是大规模预训练模型的简称，也称预训练模型或基础模型。所谓预训练模型，是指在一个原始任务上预先训练出一个初始模型，然后在下游任务中对该模型进行精调，以提高下游任务的准确性。大规模预训练模型则是指模型参数达到千亿、万亿级别的预训练模型。此类大模型因具备更强的泛化能力，能够沉淀行业经验，并更高效、准确地获取信息。
大模型的计量单位 token 指的是什么	令牌（Token）是指模型处理和生成文本的基本单位。token 可以是词或者字符的片段。模型的输入和输出的文本都会被转换成 token，然后根据模型的概率分布进行采样或计算。 例如，在英文中，有些组合单词会根据语义拆分，如 overweight 会被设计为 2 个 token：“over”、“weight”。在中文中，有些汉字会根据语义被整合，如“等于”、“王者荣耀”。 在盘古大模型中，以 N1 系列模型为例，盘古 1 token ≈ 0.75 个英文单词，1 token ≈ 1.5 汉字。不同模型的具体情况详见表 6-1。

表 6-1 token 比

模型规格	token 比（token/英文单词）	token 比（token/汉字）
N1 系列模型	0.75	1.5
N2 系列模型	0.88	1.24
N4 系列模型	0.75	1.5

训练相关概念

表 6-2 训练相关概念说明

概念名	说明
自监督学习	自监督学习 (Self-Supervised Learning, 简称SSL) 是一种机器学习方法, 它从未标记的数据中提取监督信号, 属于无监督学习的一个子集。该方法通过创建“预设任务”让模型从数据中学习, 从而生成有用的表示, 可用于后续任务。它无需额外的人工标签数据, 因为监督信号直接从数据本身派生。
有监督学习	有监督学习是机器学习任务的一种。它从有标记的训练数据中推导出预测函数。有标记的训练数据是指每个训练实例都包括输入和期望的输出。
LoRA	局部微调 (LoRA) 是一种优化技术, 用于在深度学习模型的微调过程中, 只对模型的一部分参数进行更新, 而不是对所有参数进行更新。这种方法可以显著减少微调所需的计算资源和时间, 同时保持或接近模型的最佳性能。
过拟合	过拟合是指为了得到一致假设而使假设变得过度严格, 会导致模型产生“以偏概全”的现象, 导致模型泛化效果变差。
欠拟合	欠拟合是指模型拟合程度不高, 数据距离拟合曲线较远, 或指模型没有很好地捕捉到数据特征, 不能够很好地拟合数据。
损失函数	损失函数 (Loss Function) 是用来度量模型的预测值 $f(x)$ 与真实值 Y 的差异程度的运算函数。它是一个非负实值函数, 通常使用 $L(Y, f(x))$ 来表示, 损失函数越小, 模型的鲁棒性就越好。

推理相关概念

表 6-3 训练相关概念说明

概念名	说明
温度系数	温度系数 (temperature) 控制生成语言模型中生成文本的随机性和创造性, 调整模型的softmax输出层中预测词的概率。其值越大, 则预测词的概率的方差减小, 即很多词被选择的可能性增大, 利于文本多样化。
多样性与一致性	多样性和一致性是评估LLM生成语言的两个重要方面。多样性指模型生成的不同输出之间的差异。一致性指相同输入对应的不同输出之间的一致性。
重复惩罚	重复惩罚 (repetition_penalty) 是在模型训练或生成过程中加入的惩罚项, 旨在减少重复生成的可能性。通过在计算损失函数 (用于优化模型的指标) 时增加对重复输出的惩罚来实现的。如果模型生成了重复的文本, 它的损失会增加, 从而鼓励模型寻找更多样化的输出。

提示词工程相关概念

表 6-4 提示词工程相关概念说明

概念名	说明
提示词	提示词（Prompt）是一种用于与AI人工智能模型交互的语言，用于指示模型生成所需的内容。
思维链	思维链（Chain-of-Thought）是一种模拟人类解决问题的方法，通过一系列自然语言形式的推理过程，从输入问题开始，逐步推导至最终输出结论。
Self-instruct	Self-instruct是一种将预训练语言模型与指令对齐的方法，允许模型自主生成数据，而不需要大量的人工标注。

7 安全

7.1 责任共担

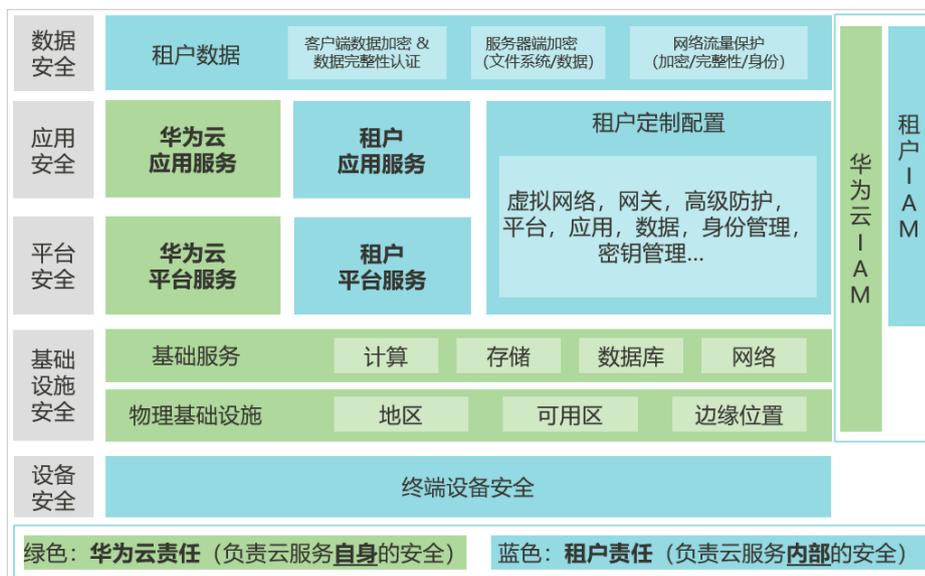
华为云秉承“将对网络和业务安全性保障的责任置于公司的商业利益之上”。针对层出不穷的云安全挑战和无孔不入的云安全威胁与攻击，华为云在遵从法律法规业界标准的基础上，以安全生态圈为护城河，依托华为独有的软硬件优势，构建面向不同区域和行业的完善云服务安全保障体系。

安全性是华为云与您的共同责任，如[图7-1](#)所示。

- **华为云**：负责云服务**自身**的安全，提供安全的云。华为云的安全责任在于保障其所提供的IaaS、PaaS和SaaS类云服务自身的安全，涵盖华为云数据中心的物理环境设施和运行其上的基础服务、平台服务、应用服务等。这不仅包括华为云基础设施和各项云服务技术的安全功能和性能本身，也包括运维运营安全，以及更广义的安全合规遵从。
- **租户**：负责云服务**内部**的安全，安全地使用云。华为云租户的安全责任在于对使用的IaaS、PaaS和SaaS类云服务内部的安全以及对租户定制配置进行安全有效的管理，包括但不限于虚拟网络、虚拟主机和访客虚拟机的操作系统，虚拟防火墙、API网关和高级安全服务，各项云服务，租户数据，以及身份账号和密钥管理等方面的安全配置。

《[华为云安全白皮书](#)》详细介绍华为云安全性的构建思路与措施，包括云安全战略、责任共担模型、合规与隐私、安全组织与人员、基础设施安全、租户服务与租户安全、工程安全、运维运营安全、生态安全。

图 7-1 华为云安全责任共担模型



7.2 身份认证与访问控制

用户可以通过调用REST网络的API来访问盘古大模型服务，有以下两种调用方式：

- Token认证：通过Token认证调用请求。
- AK/SK认证：通过AK（Access Key ID）/SK（Secret Access Key）加密调用请求。经过认证的请求总是需要包含一个签名值，该签名值以请求者的访问密钥（AK/SK）作为加密因子，结合请求体携带的特定信息计算而成。通过访问密钥（AK/SK）认证方式进行认证鉴权，即使用Access Key ID（AK）/Secret Access Key（SK）加密的方法来验证某个请求发送者身份。

7.3 数据保护技术

盘古大模型服务通过多种数据保护手段和特性，保障存储在服务中的数据安全可靠。

表 7-1 盘古大模型的数据保护手段和特性

数据保护手段	简要说明
传输加密（HTTPS）	盘古服务使用HTTPS传输协议保证数据传输的安全性。
基于OBS提供的数据保护	基于OBS服务对用户的数据进行存储和保护。请参考OBS数据保护技术说明： https://support.huaweicloud.com/productdesc-obs/obs_03_0375.html

7.4 审计

云审计服务（Cloud Trace Service, CTS）是华为云安全解决方案中专业的日志审计服务，提供对各种云资源操作记录的收集、存储和查询功能，可用于支撑安全分析、合规审计、资源跟踪和问题定位等常见应用场景。

用户开通云审计服务并创建、配置追踪器后，CTS可记录用户使用盘古的管理事件和数据事件用于审计。

CTS的详细介绍和开通配置方法，请参见[CTS快速入门](#)。

8 权限管理

如果您需要对华为云上购买的盘古大模型资源，为企业中的员工设置不同的访问权限，以达到不同员工之间的权限隔离，您可以使用统一身份认证服务（IAM）和盘古角色管理功能进行精细的权限管理。

如果华为云账号已经能满足您的要求，不需要创建独立的IAM用户（子用户）进行权限管理，您可以跳过本章节，不影响您使用服务的其他功能。

通过IAM，您可以在华为云账号中给员工创建IAM用户（子用户），并授权控制他们对华为云资源的访问范围。例如，您的员工中有负责软件开发的人员，您希望他们拥有接口的调用权限，但是不希望他们拥有训练模型或者访问训练数据的权限，那么您可以先创建一个IAM用户，并设置该用户在盘古平台中的角色，控制对资源的使用范围。

IAM 权限

默认情况下，管理员创建的IAM用户（子用户）没有任何权限，需要将其加入用户组，并对用户组授权，才能使得用户组中的用户获得对应的权限。授权后，用户就可以基于被授予的权限对云服务进行操作。

服务使用OBS存储训练数据和评估数据，如果需要对OBS的访问权限进行细粒度的控制。可以在盘古服务的委托中增加Pangu OBSWriteOnly、Pangu OBSReadOnly策略，控制OBS的读写权限。

表 8-1 策略信息

策略名称	拥有细粒度权限/Action	权限描述
Pangu OBSWriteOnly	obs:object:AbortMultipartUpload obs:object:DeleteObject obs:object:DeleteObjectVersion obs:object:PutObject	拥有用户OBS桶写权限

策略名称	拥有细粒度权限/Action	权限描述
Pangu OBSReadOnly	obs:bucket:GetBucketLocation obs:bucket:HeadBucket obs:bucket:ListAllMyBuckets obs:bucket:ListBucket obs:object:GetObject obs:object:GetObjectAcl obs:object:GetObjectVersion obs:object:GetObjectVersionAcl obs:object:ListMultipartUploadParts	拥有用户OBS桶只读权限

盘古用户角色

盘古大模型的用户可以被赋予不同的角色，对平台资源进行精细化的控制。

表 8-2 角色定义

角色名称	角色描述
超级管理员	订购服务的用户，具备当前平台下对所有工作空间的所有权限。
管理员	对工作空间有完全访问权，包括查看、创建、编辑或删除（适用时）工作空间中的资产，同时拥有添加、移除所在空间成员以及编辑所在空间成员角色的权限。
模型开发工程师	可以执行模型开发工具链模块的所有操作，但是不能创建或者删除计算资源，也不能修改所在空间本身。
应用开发工程师	应用开发工程师具备执行应用开发工具链模块所有操作的权限，其余角色不具备。
标注管理员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据加工 > 标注任务 > 任务管理” 模块 “数据工程 > 数据加工 > 标注任务 > 标注作业” 模块 “数据工程 > 数据加工 > 标注任务 > 标注审核” 模块 “数据工程 > 数据管理 > 数据集” 模块
标注作业员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据加工 > 标注任务 > 标注作业” 模块

角色名称	角色描述
标注审核员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据加工 > 标注任务 > 标注审核”模块
评估管理员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据管理 > 数据集”模块 “数据工程 > 数据管理 > 数据评估 > 人工评估”模块 “数据工程 > 数据管理 > 数据评估 > 人工评估标准”模块
评估作业员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据管理 > 数据评估 > 人工评估”模块
数据导入员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据获取 > 导入任务”模块 “数据工程 > 数据管理 > 数据集”模块
数据加工员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据加工 > 加工任务”模块 “数据工程 > 数据加工 > 合成任务”模块 “数据工程 > 数据加工 > 配比任务”模块 “数据工程 > 数据管理 > 数据指令”模块 “数据工程 > 数据管理 > 数据集”模块
数据发布员	拥有权限如下： <ul style="list-style-type: none"> “数据工程 > 数据发布 > 发布任务”模块 “数据工程 > 数据管理 > 数据集”模块

9 约束与限制

本节介绍盘古大模型服务在使用过程中的约束和限制。

规格限制

盘古大模型服务的规格限制详见[表9-1](#)。

表 9-1 规格限制

资产、资源类型	规格	说明
模型资产、数据资源、训练资源、推理资源	所有按需计费、包年/包月中的模型资产、数据资源、训练资源、推理资源。	购买的所有类型的资产与资源仅支持在西南-贵阳一区域使用。

配额限制

盘古大模型服务的配额限制详见[表9-2](#)。

表 9-2 配额限制

资源类型	默认配额限制	是否支持调整
模型实例	ModelArts Studio平台上，单个用户最多可创建和管理2000个模型实例。	是 如果希望申请提升配额，请联系客服。

功能限制

盘古大模型服务的功能限制详见[表9-3](#)。

表 9-3 功能限制

功能类型	使用限制
数据工程-数据格式要求	ModelArts Studio平台支持接入的数据需要满足格式要求，包括文件格式、单个文件大小、所有文本大小以及文件数量等，请参考《用户指南》“使用数据工程构建数据集 > 数据集格式要求”。
模型开发-训练、评测最小数据量要求	使用ModelArts Studio平台训练、评测不同模型时，存在不同数据量的限制。以NLP大模型为例，请参考《用户指南》“开发盘古NLP大模型 > 使用数据工程构建NLP大模型数据集”。
模型开发-模型最小训练单元	不同模型的最小训练单元有所不同，具体信息请参见 模型能力与规格 。
模型开发-NLP大模型请求的最大Token数	不同系列的NLP大模型支持请求的最大Token数有所不同，具体信息请参见 模型能力与规格 。

10 与其他服务的关系

与对象存储服务的关系

盘古大模型使用对象存储服务（Object Storage Service，简称OBS）存储数据和模型，实现安全、高可靠和低成本存储需求。

与 ModelArts 服务的关系

盘古大模型使用ModelArts服务进行算法训练部署，帮助用户快速创建和部署模型。

与云搜索服务的关系

盘古大模型使用云搜索服务CSS，加入检索模块，提高模型回复的准确性、解决内容过期问题。