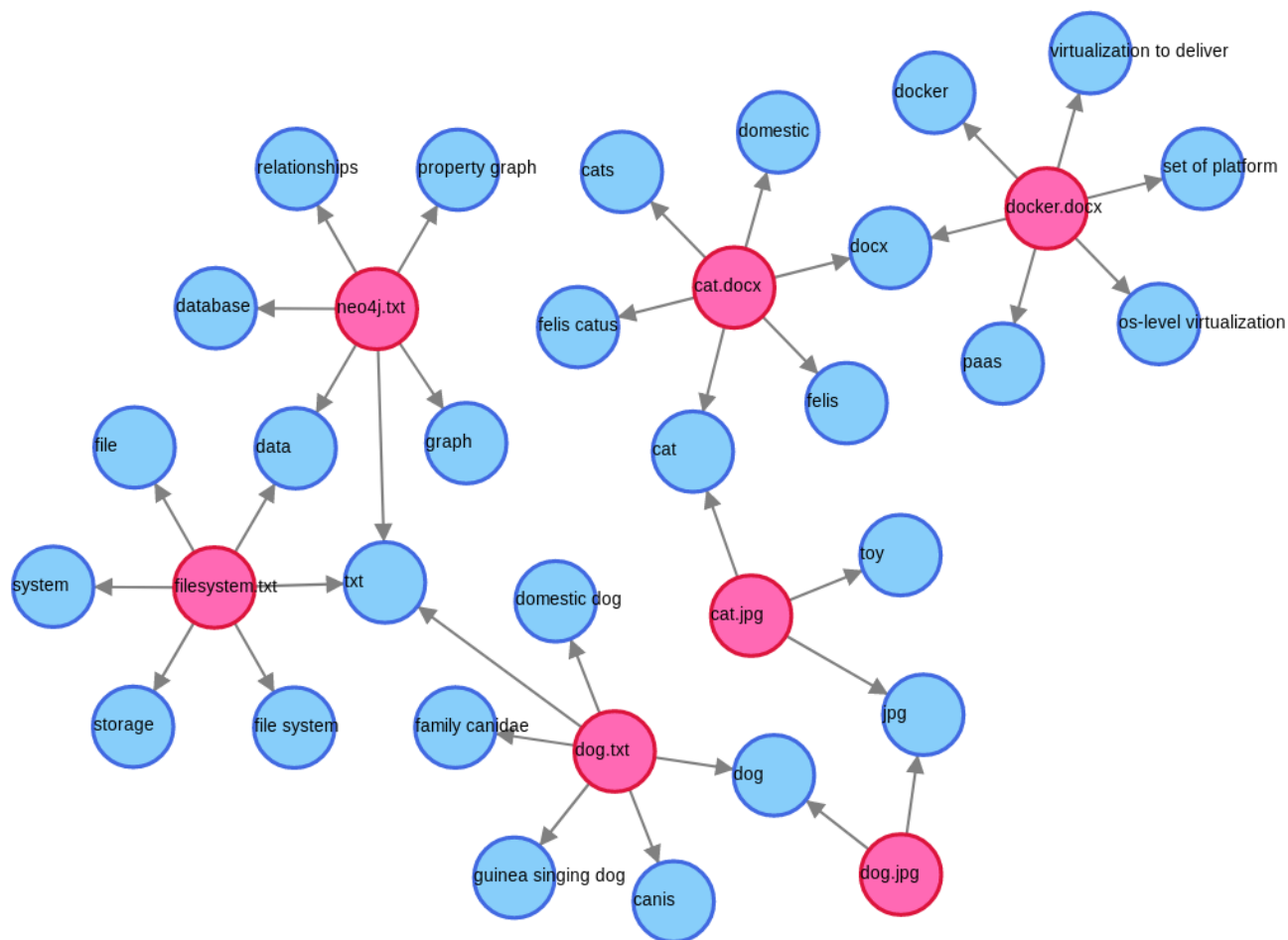


# 往届项目调研报告

## x-DisGraFS（2021）与 GDBFS（2020）

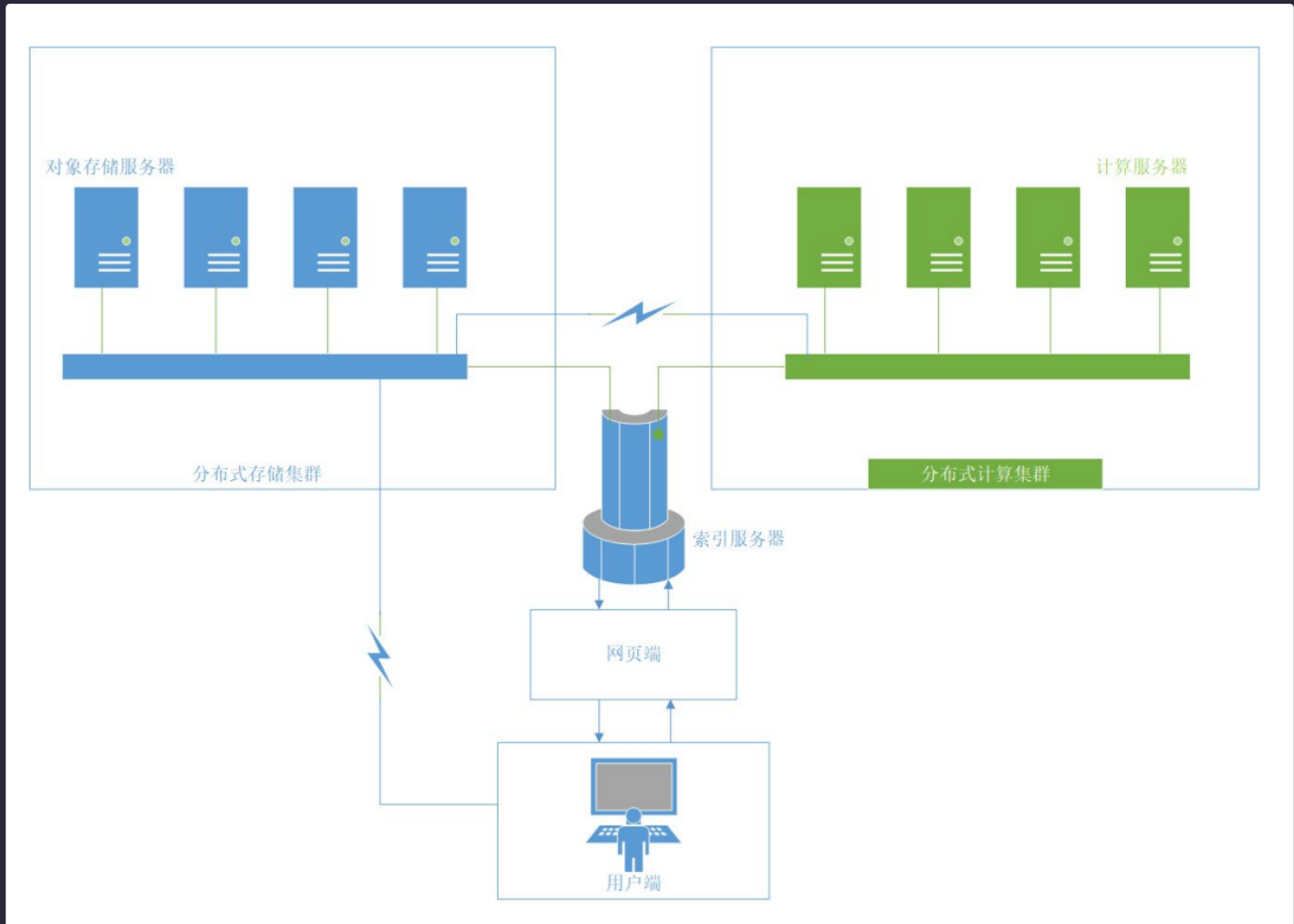
### 项目简介

GDBFS 是一个基于图数据库的文件管理系统，其设计初衷是使得人类查找文件的过程更符合人的思维——基于信息之间的相关性。



从以上图片可以看出，在图文件管理系统中，文件被自动提取出了许多属性，通过属性即可在不同文件间建立联系。

x-DisGraFS 是 GDBFS 的升级版，从单机变成了分布式集群。



大致可分为三类服务器，分布式存储服务器基于 Juicefs 管理、存储和调度分布式存储系统中的所有文件，分布式计算集群基于 Ray 将文本语义识别、语音识别等任务分散给集群中的多个计算机，索引服务器负责集群间的通信与图数据库的维护。

## 主要技术

- 图数据库

图数据库 (graph database) 是一个使用图结构进行查询的数据库，它使用节点、边和属性来表示和存储数据。该系统的关键概念是图，它直接将存储中的数据项，与数据节点和节点间表示关系的边的集合相关联。这些关系允许直接将存储区中的数据链接在一起，并且在许多情况下，可以通过一个操作进行检索。图数据库将数据之间的关系作为优先级。查询图数据库中的关系很快，因为它们永久存储在数据库本身中。可以使用图数据库直观地显示关系，使其对于高度互连的数据非常有效。

项目中使用的图数据库是 Neo4j

- 文件管理系统

文件系统提供在存储介质上组织数据的一种方式方法。其功能包括：管理和调度文件的存储空间，提供文件的逻辑结构、物理结构和存储方法；实现文件从标识到实际地址的映射，实现文件的控制操作和存取操作，实现文件信息的共享并提供可靠的文件保密和保护措施，提供文件的安全措施。

GDBFS 使用 fuse, x-DisGraFS 使用 JuiceFS

- 分布式计算

随着计算技术的发展，有些应用需要非常巨大的计算能力才能完成，如果采用集中式计算，需要耗费相当长的时间来完成。分布式计算将该应用分解成许多小的部分，分配给多台计算机进行处理。这样可以节约整体计算时间，大大提高计算效率。

Ray 是一款开源于2017年的分布式高性能计算引擎。Ray 的结构由两部分组成：application 层和 system 层 Application 层实现 API 和计算模型，执行分布式计算任务。System 层负责任务调度和数据管理，来满足表现性能和容错的要求。我们使用 Ray 来搭建分布式计算集群。

- 前端

## 难度评价

需要了解多种技术，了解分布式集群的搭建方式，难度中上

# x-ridiculous-includeos (2019)

## 项目简介

### unikernel

Unikernel是通过使用专门的库操作系统来构建的单地址空间机器镜像。开发者通过选择栈模块和一系列最小依赖库来运行应用，而这些栈和库对应于操作系统中运行应用所必需的依赖。这些库负责应用和配置代码编译，构建成封闭的、固定用途的镜像（Unikernel）可以直接在虚拟机管理程序（hypervisor）或硬件上运行，不需要类似Linux或Windows的操作系统介于其中。

---- 维基百科: [Unikernel](#)

简而言之，Unikernel 试图抹去现代操作系统带来的一些复杂性。因为“通用”的操作系统（就像任何Linux和Windows的发行版），通常会伴随着带来一些对你的应用来说并不需要的驱动、依赖包、服务、等等，但这些对每一个操作系统来说某种程度上又是必需的。

### includeOS

IncludeOS 是在云中运行 C++ 服务的 unikernel 操作系统。它提供了一个引导加载程序、标准库以及运行服务的构建和部署系统。在 VirtualBox 或 QEMU 中进行测试，并在 OpenStack 上部署服务。

简而言之，includeos 是 unikernel 的一种基于 cpp 的实现，因为 Unikernel 又可以称为 LibraryOS，它是与某种语言紧密相关的，一种 unikernel 只能用一种语言写程序，这个 LibraryOS 加上你自己写的程序最终被编译成一个操作系统。

## IoT

IoT(Internet of Things)即物联网，它的定义为把所有物品通过射频识别等信息传感设备与互联网连接起来，实现智能化识别和管理。物联网通过智能感知、识别技术与普适计算、泛在网络的融合应用，被称为继计算机、互联网之后世界信息产业发展的第三次浪潮。物联网被视为互联网的应用拓展，应用创新是物联网发展的核心，以用户体验为核心的创新2.0是物联网发展的灵魂。自2009年8月温家宝总理提出“感知中国”以来，物联网被正式列为国家五大新兴战略性新兴产业之一，写入“政府工作报告”，物联网在中国受到了全社会极大的关注。

物联网将很多领域联系到了一起，使得很多不同的领域能够共享时间或者位置等数据。如果时间与位置不正确，或是因不可靠、高延迟、或不安全的系统建置而无法正常通讯，那么，物联网的许多很重要的功能或者优势便无法发挥效果。不准确的位置以及150微秒或更长的延迟，对消费装置和家庭网络或许可以被接受，但是对于工业应用，从效能、耐用度、安全性以及可靠度等方面来考虑，是不可接受的。

因此，物联网的发展需要更强的及时性，如何降低延迟，增强通信的实时性和准确性是工程师们一直在研究的重要课题。

## 工作内容

高度概括一下就是利用 `inlcudeOS` 的低延迟性解决 `IoT` 设备中延迟高的问题，然而现今 `includeOS` 不支持 `arm` 架构，故该项目所做的工作便是构建基于 `arm` 的 `includeOS`

## 主要技术

- `unikernel`
- `includeos`
  - 读源码
- `Aarch64`体系结构
  - `mmu`、`assembly`、`system registers`、`interrupt`、`device...`
- 代码构建工具
  - `makefile`、`cmake...`

## 难度评价

极高，劝退

## 个人设计思路

偏应用层 + 分布式

- `elasticsearch`

往届项目的设计思路一般是从某个实际问题出发寻求解决办法，我们可以在[美团技术团队](#)上看到许多有关美团在前沿技术上的工作与未来展望，或许会对我们的设计有帮助

# 未来展望：构建云原生操作系统

我们认为，云原生时代的集群管理，会从之前的管理硬件、资源等职能全面转变为以应用为中心的云原生操作系统。以此为目标，美团集群调度系统还需从以下几方面发力：

- **应用链路交付管理**。随着业务规模和链路复杂度的增大，业务所依赖的PaaS组件和底层基础设施的运维复杂度早已超过普遍认知，对于刚接手项目的新人更是难上加难。所以我们需要支持业务通过声明式配置交付服务并实现自运维，给业务提供更好的运维体验，提升应用的可用性和可观测性，减少业务对底层资源管理的负担。
- **边缘计算解决方案**。随着美团业务场景的不断丰富，业务对边缘计算节点的需求增长，比预期快很多。我们会参考业内最佳实践，形成适合在美团落地的边缘解决方案，尽快为有需求的服务提供边缘计算节点管理能力，实现云边端协同。
- **在离线混部能力建设**。在线业务集群的资源利用率提升是有上限的，根据Google在论文《Borg: the Next Generation》中披露的2019年数据中心集群数据，刨去离线任务，在线任务的资源利用率仅为30%左右，这也说明了再往上提升风险较大，投入产出比不高。后续，美团集群调度系统将持续探索在离线混部，不过由于美团的离线机房相对独立，我们的实施路径会与业界的普遍方案有所不同，会先从在线服务和近实时任务的混部开始，完成底层能力的构建，再探索在线任务和离线任务的混部。

## 5 总结与展望

TensorFlow在大规模推荐系统中被广泛使用，但由于缺乏大规模稀疏的大规模分布式训练能力，阻碍了业务的发展。美团基于TensorFlow原生架构，支持了大规模稀疏能力，并从多个角度进行了深度优化，做到千亿参数、千亿样本高效的分布式训练，并在美团内部进行了大规模的使用。对于这类关键能力的缺失，TensorFlow社区也引起了共鸣，社区官方在2020年创建了SIG Recommenders[11]，通过社区共建的方式来解决此类问题，美团后续也会积极的参与到社区的贡献当中去。

美团推荐系统场景的模型训练，目前主要运行在CPU上，但随着业务的发展，有些模型变得越来越复杂，CPU上已经很难有优化空间（优化后的Worker CPU使用率在90%以上）。而近几年，GPU的计算能力突飞猛进，新一代的NVIDIA A100 GPU，算力达到了156TFLOPS（TF32 Tensor Cores）、80G显存、卡间带宽600GB/s。对于这类复杂模型的Workload，我们基于A100 GPU架构，设计了下一代的分布式训练架构，经过初步优化，在美团某大流量业务推荐模型上也拿到了较好的效果，目前还在进一步优化当中，后续我们会进行分享，敬请期待。