

EDUCATION

- **Master - EE - Beihang University** Beijing, PRC
Exempt Entry Exam; Finish courses and research publishing mission during the 1st semester. Sep. 2016 – Jan. 2019
- **Bachelor - EE - Beihang University** Beijing, PRC
Honors College(top 50 among 3000+ freshmen.); GPA 3.7. Sept. 2012 – Jun. 2016
- **High School - Dalian Yuming High School** Dalian, PRC
CPhO - Provincial First Prize, granted exemption from the college entrance examination. Sept. 2012 – Jun. 2016

COMPETITIVE PROGRAMMING

- **Google Code Jam Kickstart 2017**
Top 5%, rank 108th globally. [scoreboard](#), [id: WeiYong1024](#)

CAREER

- **Mindverse(Joined angel round, a pioneering startup for LLM app)** Hangzhou, PRC
Infra Tech Lead(4-person team), building AiInfra+DataInfra+DevOps for LLM applications. Jun. 2022 - Now
 - **Company representative products**
 - * **MeBot**: Inspiring personal assistant, Web version ranked second in the ProductHunt weekly list in August 2024, App version is now available on the Apple App Store.
 - * **MindOS**: Early LLMAgentPlatform, launched in October 2022, ranked first in ProductHunt's weekly list in July 2023.
 - **Company technical design**
 - * **LPM**: Large personal model. Based on Lora and user online data to train personalized LLM, which can be used in production environments.
 - * **MindOS**: The earliest LLM Agent Platform in China, launched 17 months earlier than ByteDance's Coze platform.
 - * **UMM**: Unified mind model. 16 months earlier than OpenAI first proposed the concept of AIAgent.
 - **Infra Technology Responsibilities**
 - * **AiInfra - Best practices for building infra of LLM based application**
 - **MLOps Platform**: Design and implement a full-link automated MLOps platform that includes data collection, data augmentation, LLM training, and LLM deployment. Combine the model observation capabilities of the open-source tool ClearML to support real-time LPM training and production deployment..
 - **LLM Unified Access Layer**: In the early stage, Java was used to build basic services, achieving unified access, deployment management, multi-account pooling, health check, weight allocation, cost statistics, and visualization for models such as OpenAI, AzureOpenAI, GCPVertexAI, and Claude. In the later stage, it was migrated to LiteLLM as the best practice to provide the above functions.
 - **LLM Engineering Platform**: Use Langfuse's standardized algorithm for offline experiment management, prompt asset management, asset version management, and online LLM-Tracing functions.
 - * **DevOps - Efficient and safe production**
 - **Production Environment Management**: Designed and built multiple production environments such as Prod/Pre/Test, including scalable GitServer with strict CodeReview processes, CI/CD processes, zero-trust device management, and other infrastructure. Technologies involved: Gitlab, Helm, Java, Jenkins, JumpServer, Kubernetes, Octant, Python, Rancher, Shell, etc.
 - **System Observability**: Selection, design, and construction of the distributed configuration, distributed rate limiting, distributed scheduling, application performance monitoring, and log monitoring internal production infrastructure used by the company's backend system. Technologies involved: Apollo, Datadog, Grafana, Loki, Nacos, Prometheus, Sentinel, Skywalking, xxl-job, Tencent Cloud CLS, etc.

- **Cloud and third-party service management:** Internal user permission management, refined control of cloud usage and third-party service costs, design, production, and maintenance cost dashboard. Design technologies: Azure, GCP, CronJob, TencentCloud, etc.
- * **DataInfra - Data assets, experimental tools, data-driven decision making**
 - **Online and Offline Data System:** Design, build, and maintain online and offline data streams including OLTP->CDC->OLAP->DataVisualization for offline data development, and asynchronous data streams including message queue->buried point service for user behavior data analysis and other scenarios. Support various metrics dashboards, operational BI reports, and other requirements. Technologies involved: MySQL, Clickhouse, Debezium, Superset, RocketMQ, etc.
 - **A/B Testing Tools:** Selection, design, construction, and maintenance of a full-link experimental tool system that includes experimental configuration/result analysis platform (GrowthBook) + front-end tracking (Google Analytics) + back-end tracking (Springboot implementation, including engineering and algorithms). Implement feature toggles and A/B testing for front-end pages and back-end algorithms. Technologies involved: Google Analytics, GrowthBook, Java, NodeJS, Python, etc.
 - **Product Privacy Compliance:** Establish a series of privacy compliance standards, including but not limited to data anonymization, internal training, etc., to ensure that the company's products meet the privacy compliance standards of North America's USDP and the EU's GDPR.
- * **Business Layer Development**
 - **Middle Platform Services:** Designed and built mid-tier basic microservices, offering common foundational services to business systems like backend encrypted storage, user privacy database, web crawlers, offline data warehouse management, etc., and delivered them in forms of RestfulAPI, secondary libraries + RPC interfaces, etc.

○ Business Support

- * Led mobile app mini-programs and web product MindOS's subsystems & full-link load testing, and safeguarded the global launch process; prepared and established SOPs for toB public cloud/private cloud/private deployment; supported service architecture migration across countries/cloud regions as per business needs.
- * Promoted a company-wide blameless postmortem culture, organized safety production weekly meetings.

• **Alibaba Cloud**(China's leading cloud service provider)

Senior Engineer, Backend, Java, Springboot, Alibaba Cloud

Hangzhou, PRC

Aug. 2020 - May. 2022

- **Traffic Scheduling Middleware:** : In the group's high-availability team, responsible for Alibaba Cloud's internal traffic scheduling middleware. Utilizing Alibaba's internal PandoraBoot (an enhanced version of Springboot), provided container granularity second-level health monitoring and 3 σ outlier detection on the cluster. Leveraged the weight table mechanism of Alibaba's internal RPC framework HSF for minute-level traffic scheduling and recovery linkage of microservices. This system prevented cluster avalanches caused by single-point failures and was also used for JIT preheating. During my tenure, the product was integrated into all core applications of Alibaba Group, supporting over 400,000 containers, and served as a core security component safeguarding the peak traffic of 600,000 QPS during the 2020 Double Eleven shopping festival.
- **Private Cloud Management Platform:** As part of the P8 team, the goal was to build a private cloud management system for government scenarios. Developed using the Java framework Springboot, the system is divided into three subsystems: application operations, resource operations, and general management. I was responsible for backend development of the application operations center subsystem. The subsystem's base layer included cloud resource listing, monitoring and inspection, usage statistics, custom interface inspection, SQL inspection, etc. On top of this, we built expert reports such as slow SQL statistical analysis, business dashboards, and health check reports. This system is currently commercially deployed in over 6 locations nationwide, supporting projects like elementary school admissions in a provincial capital city in 2021 and nucleic acid testing during the 2022 Spring Festival.
- **Lightweight Delivery for Government Cloud Management:** As a new P9 team member, the goal was to deliver private clouds for government use in 7 regions. Initially, developers used the Rainbond interface to deploy microservices from different teams onto customer K8S clusters, but the complicated configuration led to long delivery cycles. To address this, I developed a lightweight deployment delivery solution using docker-compose, based on generic environment initialization scripts and delivery step documents. By maintaining environment variables for different customer sites, we centralized control and integrated over 20 microservices from various older teams. I edited the WBS and SOP, organized training for outsourced staff on the delivery process, and trained outsourced individuals could deliver or upgrade a site in 2 days. With this solution, the team's monthly commercial output capacity increased from single to double digits.

- **Pony.ai(L4 autonomous driving solution)**

Engineer, Infra, C++, Python, Linux, Bazel, Strict Code Quality Standards

Beijing, PRC

Feb. 2019 - Mar. 2020

- **Driving Recording Toolchain:** Autonomous vehicle road tests required an accompanying engineer to record issues via an external keyboard. To eliminate the need for this, a driving recording toolchain was developed using microphones and physical buttons as the hardware solution. Firstly, an input device interface library evdev, licensed under BSD, was integrated into the Bazel main project. Based on this library and Linux's ALSA audio driver, a daemon was written responsible for hardware detection, listening for signals from external buttons and the vehicle system process's pipeline, triggering and stopping microphone recording, and storing issue information as audio files on the industrial computer during the drive. In the data processing stage, a speech-to-text service was set up on GCP. Python scripts were used to extract the contents of the audio files, and Google's Protobuf tool was used to serialize the original information into the QA data stream.
- **Onboard Car Voice System:** The original onboard voice module used Google's Pico TTS library to play hardcoded voice content text in real-time during driving, leading to issues of language monotony and repetitive consumption of computing resources for TTS processes. The new onboard voice system uses an audio file-based workflow, aimed at reducing the computational load on the vehicle's industrial computer and supporting voice I18N. For this, 1. A Python-based internal CLI tool, `car_sound_utils` (referred to as CLI), was developed for engineers to create and modify existing corpora and manage voice pack versions. First, an entire TTS service built on AWS using CloudFormation and Lambda functions was encapsulated in the CLI for creating audio commands. Then, commands for uploading and downloading voice files and uploading voice packs to the internal storage server were added to the CLI, along with a caching mechanism to speed up uploads and downloads. 2. In the onboard system, the main process initializes the voice module using Linux's ALSA driver to load the audio stream into RAM, and plays it when the voice module receives a playback request message.

INTERNSHIP

- **Airbnb:** In the summer of 2018, Airbnb expanded into the Chinese market, and as one of the first 8 interns in the Mainland China office, developed the domestic version of the annual host review page with full-stack tech.
- **Megvii(Face++):** During the cooperation between Megvii and VIVO in early 2018, participated in the development of the facial recognition module for the X21 model and conducted mobile model search optimization based on InceptionV3(a type of CNN) for mobile devices.

SYSTEM EXPERIENCE & TECH STACK

- **System experience:** MLOps, LLM applications, Prompt engineering, large-scale distributed systems, high availability, observability, autonomous driving, etc.
- **Tech stack:** Java, Python, Kubernetes, Shell, C++, mainstream clouds services, mainstream DevOps tools, etc.

教育背景

- 硕士 - 北京航空航天大学 - 电子与通信工程
保研; 研一上学期发核心期刊达毕业标准; 北京市优干
中国, 北京
2016 年 9 月 - 2019 年 1 月
- 本科 - 北京航空航天大学 - 电子信息工程
沈元荣誉学院 (入学 top50/3000+); 北京市三好; GPA 3.7
中国, 北京
2012 年 9 月 - 2016 年 6 月
- 高中 - 大连育明高中
CPhO (全国中学生物理竞赛) 省一等奖获高考免试保送资格
中国, 大连
2009 年 9 月 - 2012 年 6 月

编程竞赛

- Google Code Jam Kickstart (谷歌 2017 全球校招赛)
全球 108th, 中国 18th, 前 5% [计分板链接 \(id - WeiYong1024\)](#)

工作经历

- 心识宇宙 (天使轮加入, 国内 LLM 应用先驱创业团队)
基础设施技术主管, 领导 4 人团队, 构建大模型应用的 AiInfra+DataInfra+DevOps 体系
中国, 杭州
2022 年 6 月 - 至今
 - 公司代表产品
 - * MeBot: 启发型个人助理, Web 版于 2024.8 获 ProductHunt 周榜第二, App 版已上苹果 AppStore。
 - * MindOS: 国内早期 LLM Agent Platform, 2022.10 上线, 2023.7 获 ProductHunt 周榜第一。
 - 公司技术设计
 - * LPM(Large personal model): 基于 Lora 与用户在线数据训练个性化 LLM, 可用于生产环境。
 - * MindOS: 国内最早的 LLM Agent Platform, 比字节 Coze 平台上线时间早 17 个月。
 - * UMM(Unified mind model): 统一心识模型, 比 OpenAI 首次提出 AIAgent 概念早 16 个月。
 - Infra 技术职责
 - * AiInfra ——打造 LLM 应用基础设施最佳实践
 - MLOps 平台: 设计并实施包含数据采集、数据增强、LLM 训练、LLM 部署的全链路自动化 MLOps 平台。结合开源工具 ClearML 的模型观测能力, 用于支撑 LPM 实时训练并做生产部署。
 - LLM 统一接入层: 前期使用 Java 自建基础服务, 实现 OpenAI、AzureOpenAI、GCPVertexAI、Claude 等模型的统一接入、部署管理、多账号池化、探活、权重分配、成本统计与可视化等功能。后期迁移至 LiteLLM 作为提供上述功能的最佳实践。
 - LLM 工程平台: 使用 Langfuse 标准化算法离线实验管理、Prompt 资产管理、资产版本管理、线上 LLM-Tracing 等功能。
 - * DevOps ——高效、安全生产
 - 生产环境管理: 设计并搭建 Prod/Pre/Test 等多套生产环境, 包含严格 CodeReview 流程的可伸缩 GitServer、CI/CD 流程、零信任设备管理等基础设施。涉及技术: Gitlab、Helm、Java、Jenkins、JumpServer、Kubernetes、Octant、Python、Rancher、Shell 等。
 - 系统可观测性: 选型、设计并构建公司后台系统使用的分布式配置、分布式限流、分布式调度、应用性能监控及日志监控内部生产等基础设施。涉及技术: Apollo、Datadog、Grafana、Loki、Nacos、Prometheus、Sentinel、Skywalking、xxl-job、腾讯云 CLS 等。
 - 云及三方服务管理: 公司内部用户权限管理, 精细化管控用云、用三方服务的成本, 设计、制作及维护成本大盘。涉及技术: Azure、GCP、CronJob、TencentCloud 等。
 - * DataInfra ——管理数据资产、提供实验工具、数据驱动决策
 - 在离线数据体系: 设计、构建并维护包含 OLTP->CDC->OLAP->DataVisualization 的在离线数据流用于离线数据开发, 包含消息队列->埋点服务的异步数据流用于用户行为数据分析等场景。支撑各类指标的大盘、运营 BI 报表等需求。涉及技术: MySQL、Clickhouse、Debezium、Superset、RocketMQ 等。
 - A/BTest 工具体系: 选型、设计、构建并维护包含实验配置/结果分析平台 (GrowthBook) + 前端埋点 (GoogleAnalytics) + 后端埋点 (Springboot 实现, 包含工程和算法) 的全链路实验工具体系。实现前端页面、后端算法的功能开关和 A/B 实验。涉及技术: GoogleAnalytics、GrowthBook、Java、NodeJS、Python 等。

- **产品隐私合规**: 构建一系列隐私合规标准, 包括但不限于数据脱敏、内部培训等, 使公司产品达到北美 USDP 和欧盟 GDPR 隐私合规标准。

- * **业务层研发——提供业务开发依赖的基础服务**

- **中台基础服务**: 设计并构建中台基础微服务, 为业务系统提供后端加密存储、用户隐私数据库、Web 爬虫、离线数据仓库管理等通用基础服务, 并以 RestfulAPI、二方库 +RPC 接口等形式输出。

- **支撑业务**

- * 主导移动端小程序万物总动员、Web 端产品 MindOS、App MeBot 及各种子系统的全链路压测、发布过程重保; 准备并制定 toB 公有云/专有云/私有化产品交付 SOP; 服务架构随业务跨国/跨云地域迁移等。
- * 在全公司横向推进 blameless postmortem 文化, 组织安全生产周会。

- **阿里云**

高级工程师、后端、Java、Springboot、阿里云

中国, 杭州

2020 年 8 月 - 2022 年 5 月

- **流量调度中间件**: 在集团高可用团队负责阿里云内部流量调度中间件。基于阿里内部的 PandoraBoot (Springboot 加强版) 提供集群上容器粒度秒级健康指标监控、 3σ 离群点检测, 依托阿里内部 RPC 框架 HSF(Dubbo 内部版) 的权重表机制实现微服务的分钟级流量调度与恢复链路。防单点故障引起的集群雪崩, 也可用于 JIT 预热。任职期间推广产品纳管了阿里集团全部核心应用, 支撑 40w+ 容器, 作为核心安全组件护航 2020 双十一 60wQPS 流量洪峰。
- **私有云管理平台**: 时年 P8 团队目标建设政务场景下的专有云管理系统。基于 Java 框架 Springboot 开发, 该系统分应用运维、资源运营、总集管理三个子系统。我负责其中应用运维中心子系统的后端开发。应用运维中心底层构建云资源列表、监控巡检、用量统计、自定义接口巡检、SQL 巡检等原子能力, 并在上层此构建了慢 SQL 统计分析、业务大盘、体检报告等专家报表。该系统目前在全国超过 6 个局点商业化输出, 重保护航 2021 某省会城市小学入学等项目, 2022 春节核酸检测等项目。
- **政务云管轻量化交付**: 时年新 P9 团队目标交付 7 个地方政务私有云, 早期由开发同学使用 Rainbond 界面化地将来自原不同团队的微服务组合部署在客户 K8S 集群上, 配置繁琐导致交付周期长。为解决该问题我基于 docker-compose 搭建单机版轻量化部署交付方案, 基于通用的环境初始化脚本与交付步骤文档, 仅通过维护不同客户现场的环境变量以中心化管控, 纳管来自各个老团队的 20+ 个微服务。编辑 WBS 和 SOP, 组织外包同学培训交付流程, 单个经过培训的外包可以在 2 日内交付/升级一个现场。基于该方案, 团队技术产品簇的单月商业化输出能力从个位数数据点上升到两位数。

- **小马智行(L4 自动驾驶解决方案)**

工程师、Infra、C++、Python、Linux、Bazel、严苛的代码质量标准

中国, 北京

2019 年 2 月 - 2020 年 3 月

- **行车录音工具链**: 无人车路测需要一位跟车工程师通过外接键盘描述记录路测 issue, 为去掉跟车工程师构建行车录音工具链, 以麦克风和物理按钮作为硬件方案。首先在 Bazel 项目中引入基于 BSD 软件许可的输入设备接口库 evdev, 并基于该库和 Linux 的 ALSA 音频驱动编写守护进程, 负责硬件检测、监听来自外部按钮的信号和车载系统进程的管道信息、触发与停止麦克风录音, 在行车过程中将 issue 信息存储为工控机上的音频文件。在数据处理阶段, 首先在 GCP 上搭建语音转文字服务, 并用 Python 脚本将音频文件的内容提取出来, 然后使用 Google 的 Protobuf 工具将原信息序列化汇入 QA 数据流。
- **车载语音系统**: 原车载语音模块使用 Google 的 Pico TTS 库在行车过程中将硬编码的语音内容文本实时播放, 从而导致语言单一和 TTS 过程重复消耗计算资源的问题。新的车载语音系统使用基于音频文件的工作流程, 旨在减少车载工控机的计算量并支持语音 II8N。为此 1. 研发环节用 Python 编写内部 CLI 工具 car_sound_utils(以下简称 CLI), 供内部工程师创建和变更现有语料库、管理语音包版本: 首先在 CLI 中基于 AWS 上用 CloudFormation、Lambda 函数计算搭建的整套 TTS 服务封装音频创建命令, 然后在 CLI 中添加上传、下载语音文件和上传语音包到内部 storage server 的命令, 以及缓存机制加速上传下载。2. 车载系统环节, 主进程初始化语音模块时使用 Linux 的 ALSA 驱动将音频流加载进内存, 并在语音模块接到播放请求消息时播放。

实习经历

- **爱彼迎**: 2018 年夏, 爱彼迎发力中国市场, 作为大陆首批 8 个实习生, 全栈开发了当年国内版房东年度回顾页。
- **旷视科技**: 2018 年初, 旷视和 VIVO 合作期间, 参与 X21 机型人脸识别模组的研发, 做基于 InceptionV3 (一种 CNN) 移动端模型搜索优化。

系统经验与技术栈

- **系统经验**: MLOps、LLM 应用、Prompt 工程、超大规模分布式系统、高可用、可观测性、自动驾驶等
- **技术栈**: Java、Python、Kubernetes、Shell、C++、国内外主流云、主流运维工具等