

基于新型编程语言 设计实现操作系统内核

陈 渝 李明

清华大学

第五届国产嵌入式操作系统技术与产业发展论坛

2023年08月12日

个人简介

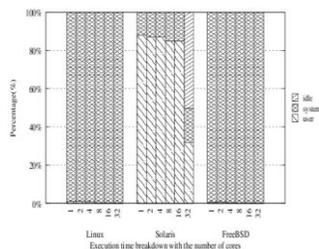


陈渝，国防科技大学本硕博毕业，清华大学长聘副教授，CSAE-CCF-CICV操作系统联合实验室首席科学家，中国汽车工程学会基础软件分会副主任委员，中国计算机学会系统软件专委会常委委员，MIT访问学者

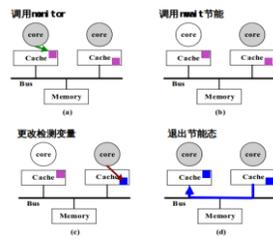
研究方向以**操作系统**为核心，涉及系统安全、性能优化、程序分析等领域。近十余年来，主要在操作系统/虚拟机性能优化，系统软件安全分析，取得了有一定国际/国内影响的研究成果。

操作系统性能优化

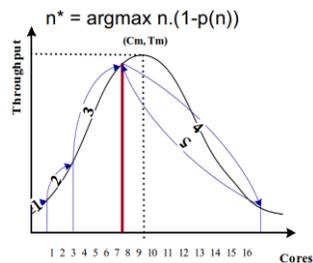
性能分析与方法



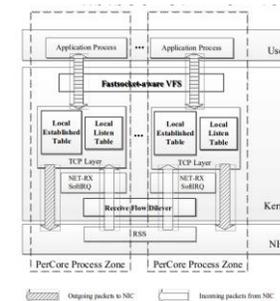
锁机制建模与可扩展优化



调度算法,内存管理

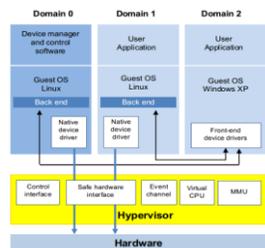


内核IO优化机制与系统

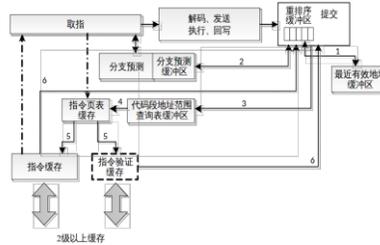


操作系统安全分析

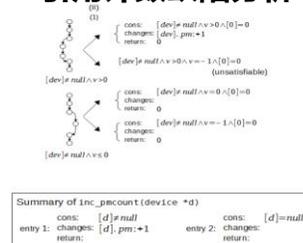
OS分区隔离安全机制



防止代码重用攻击的软硬件协同保护机制



基于不一致路径对检测的引用计数缺陷分析



目录

- 一 | 整体规划
- 二 | 科研计划
- 三 | 落地规划



研究现状

操作系统是具有战略意义的关键信息系统核心基础设施

- 操作系统是管理计算机硬件与软件资源的基础程序，是最基本、最为重要的基础性软件系统
- 被列为科技日报**35项“卡脖子”核心技术之一**（排第三）

操作系统的核心问题与应用需求

- **安全问题（重点）**：操作系统安全是计算机安全的软件基石
 - Linux等基于C语言的操作系统漏洞有扩大化的趋势
- **性能问题**：在多核、低延迟、新型硬件支持等方面需改进
 - 随着核数增多，内核锁竞争导致吞吐量下降
- **生态问题**：需要能够兼容已有生态，利于建设应用环境

面向智能设备，高安全、高性能、应用生态好的操作系统是生态（智能网联汽车）发展的必需



2023年度的世界黑客大赛，成功实现针对特斯拉的 TOCTOU 的竞态条件漏洞攻击

《丧失先机，
国之痛》
(2018年4月23日)
没有自研操作系统的大

研究现状

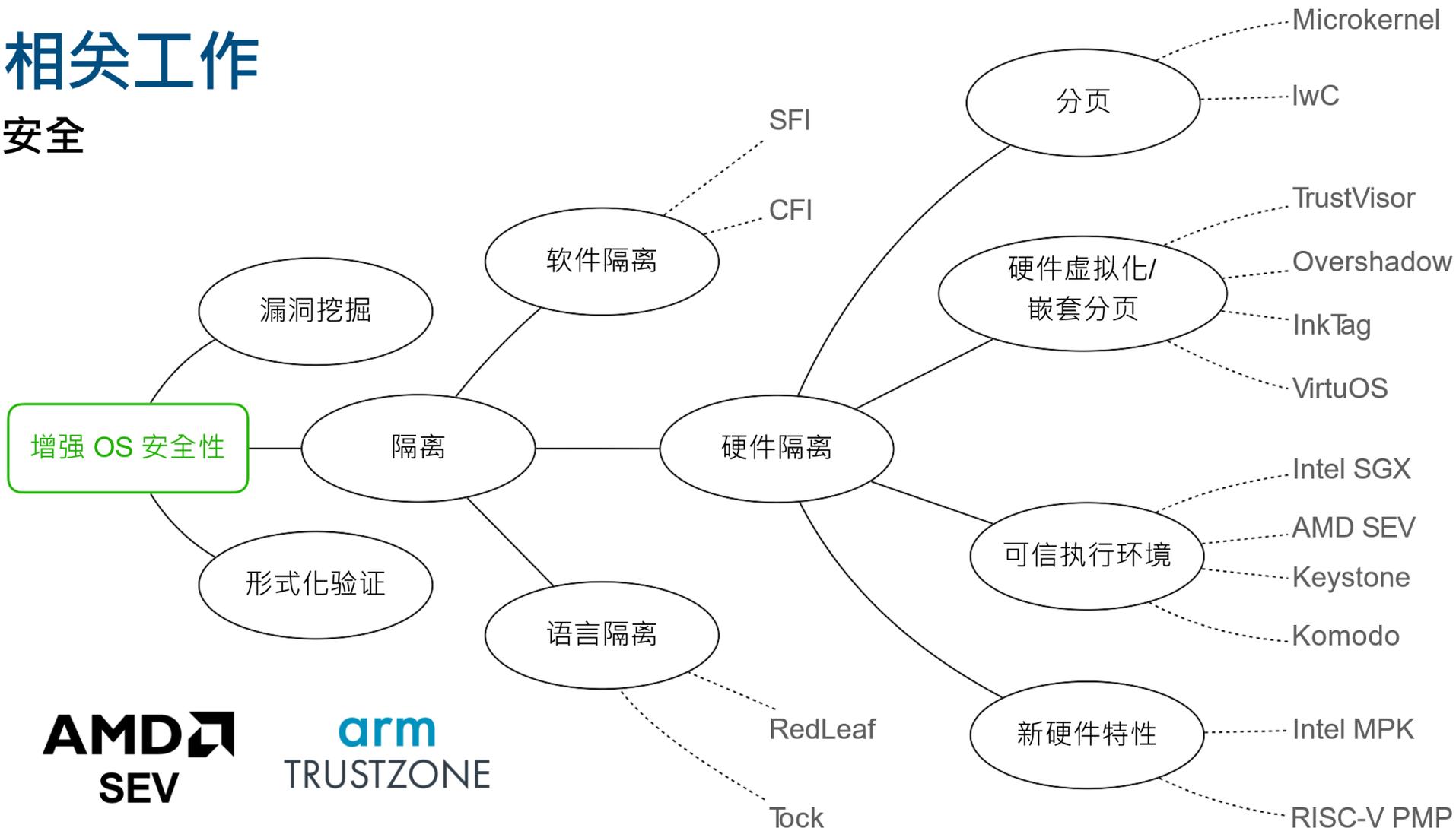
操作系统架构

编程语言

资源管理

相关工作

安全



Process-based



VM-based



Separate worlds

拟解决的关键科学问题

- 重点提出基于Rust的可扩展操作系统安全架构，较现有OS提供更多层次的安全保障能力，为构建自主可控IT生态奠定OS基础。

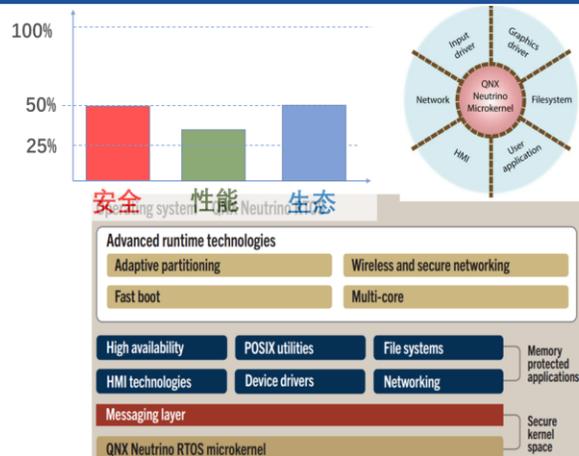
一、如何开发与设计可灵活组合、支持持续适应优化和可持续演化的操作系统安全架构？



宏内核架构 monolithic-kernel

- 伪模块化：任意两个模块都可有调用关系
- 高复杂化：模块之间存在上百个的关联性
- 历史遗留：老代码难以移除，成为包袱
- 安全漏洞：C语言是多数典型漏洞的根源
- 使用场景：移动端、桌面、服务器领域

5



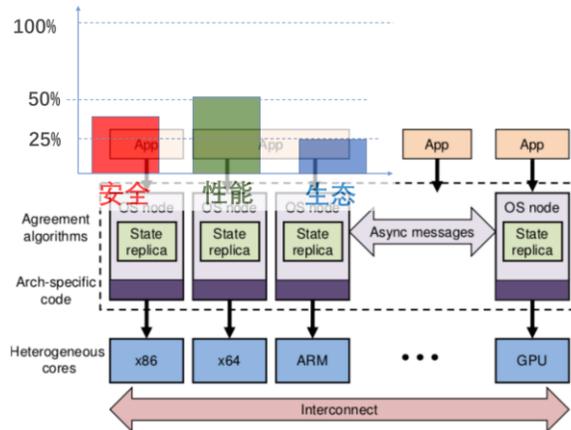
微内核架构 micro-kernel

- 应用生态：应用范围小，安全检查不足
- 高复杂化：模块之间存在上百个的关联性
- 安全漏洞：C语言是多数典型漏洞的根源
- IPC隔离：安全缺陷依然不容忽视
- 使用场景：嵌入式系统

外核架构 exo-kernel

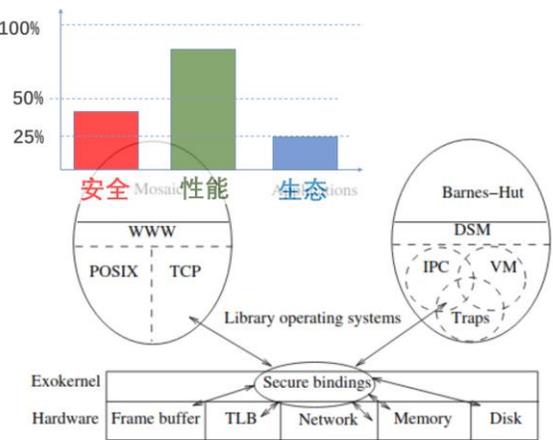
- 开发难度：不同应用需要不同的LibOS
- 安全漏洞：C语言是多数典型漏洞的根源
- 生态劣势：缺少应用生态，安全检查不足
- 使用场景：Hypervisor/VMM

7



多核架构 multi-kernel

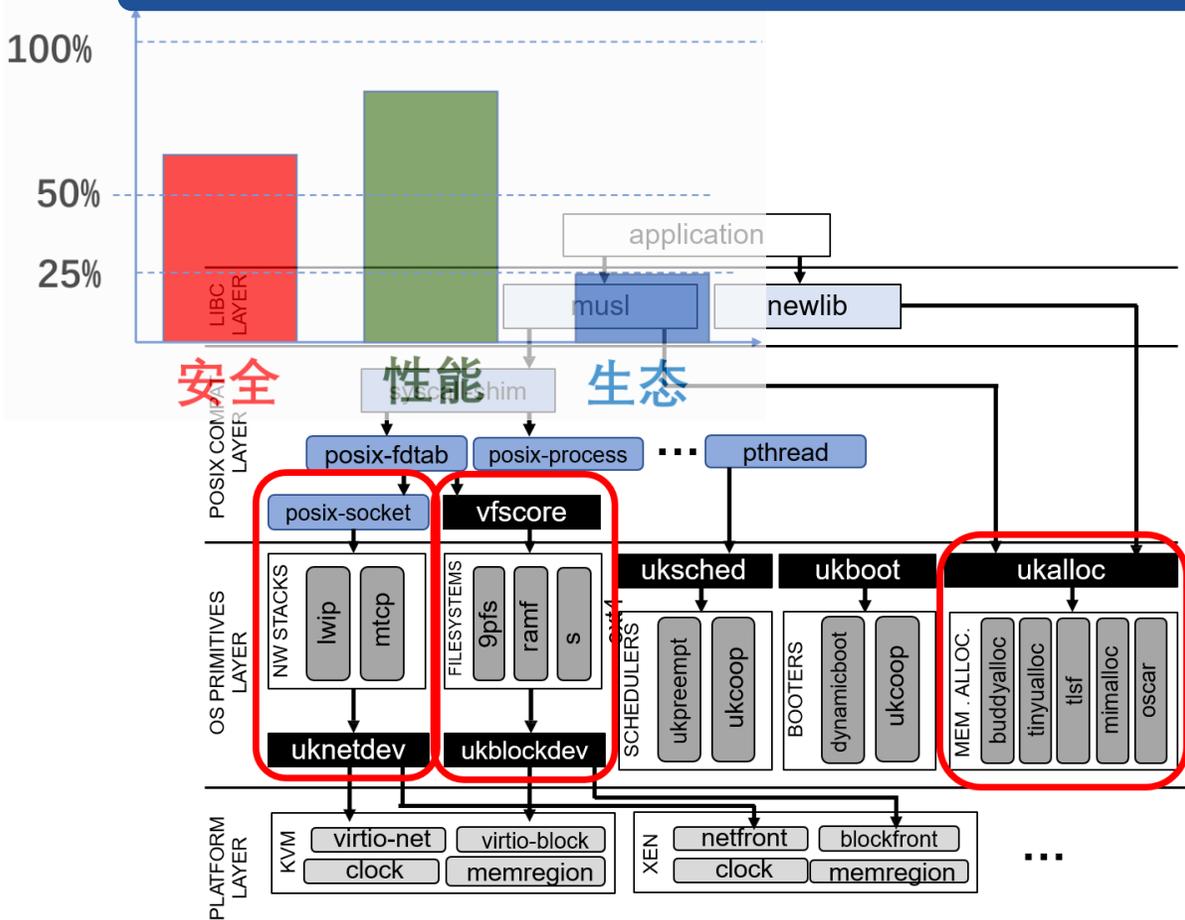
- 开发难度：分布式内核开发
- 安全漏洞：C语言是多数典型漏洞的根源
- 生态劣势：缺少应用生态，安全检查不足
- 使用场景：适合未来的异构架构



拟解决的关键科学问题

- 重点提出基于Rust的可扩展操作系统安全架构，较现有OS提供更多层次的安全保障能力，为构建自主可控IT生态奠定OS基础。

一、如何开发与设计可灵活组合、支持持续适应优化和可持续演化的操作系统安全架构？



单核架构 uni-kernel

- 开发难度：APP与内核在同一地址空间
- 安全漏洞：C语言是多数典型漏洞的根源
- 通用需求：难以通用化
- 生态劣势：缺少应用生态
- 使用场景：高性能的特定领域场景



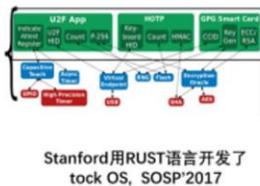
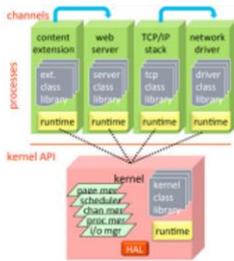
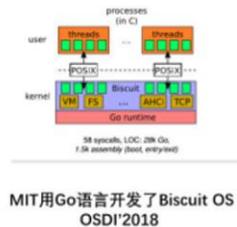
拟解决的关键科学问题

- 重点提出基于Rust的可扩展操作系统安全架构，较现有OS提供更多层次的安全保障能力，为构建自主可控IT生态奠定OS基础。

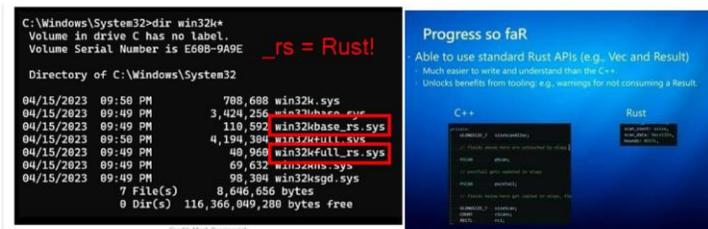
一、如何开发与设计可灵活组合、支持持续适应优化和可持续演化的操作系统安全架构？

二十一世纪以来，尝试用新语言重构操作系统的探索一直在进行，其中Ocaml, C#, Go, Rust等新一代编程语言在学术界和产业界初露头角

安全是新一代系统语言试图代替C语言的重要因素

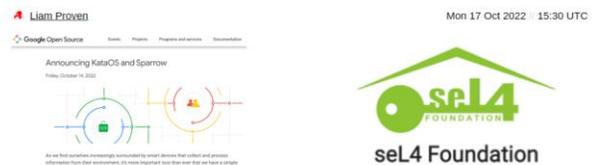


Replacing C++! Microsoft's Win11 kernel rewritten in Rust is here!



Google reveals another experimental operating system: KataOS

Based on Rust, on top of seL4 – a big deal in the microkernel world



Linus Torvalds: Rust will go into Linux 6.1

At the Kernel Maintainers Summit, the question wasn't, "Would Rust make it into Linux?" Instead, it was, "What to do about its compilers?"

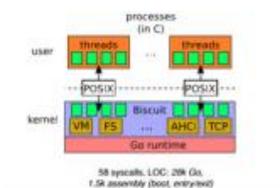


Written by Steven Vaughan-Nichols, Senior Contributing Editor on Sept. 19, 2022

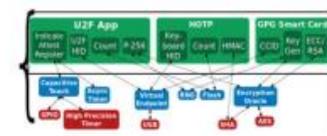
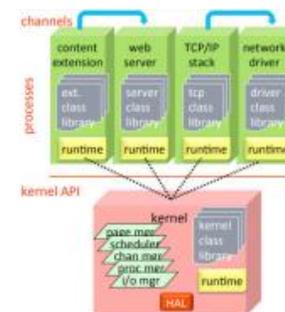


目录

- 一 | 整体规划
- 二 | 科研计划
- 三 | 落地规划



MIT用Go语言开发了Biscuit OS
OSDI'2018



Stanford用RUST语言开发了
tock OS, SOSP'2017



主要研究内容及技术架构

■理论基础:

扩展Rust语言的**所有权(Ownership)**机制到操作系统资源管理中, 原理基于:

- 类型系统的Affine types(仿射类型)
- 值/资源拥有者唯一/所有权转移/借用

核心规则:

- **所有权规则**: 每个值/资源在任意时刻只能由一个所有者拥有
- **借用规则**: 临时地借用值而不获取所有权
- **生命周期规则**: 验证引用或借用的有效范围

基于Affine Type的Ownership安全机制

Rust核心库



	Exchange	Weakening	Contraction	Use
Ordered	—	—	—	Exactly once in order
Linear	Allowed	—	—	Exactly once
Affine	Allowed	Allowed	—	At most once
Relevant	Allowed	—	Allowed	At least once
Normal	Allowed	Allowed	Allowed	Arbitrarily

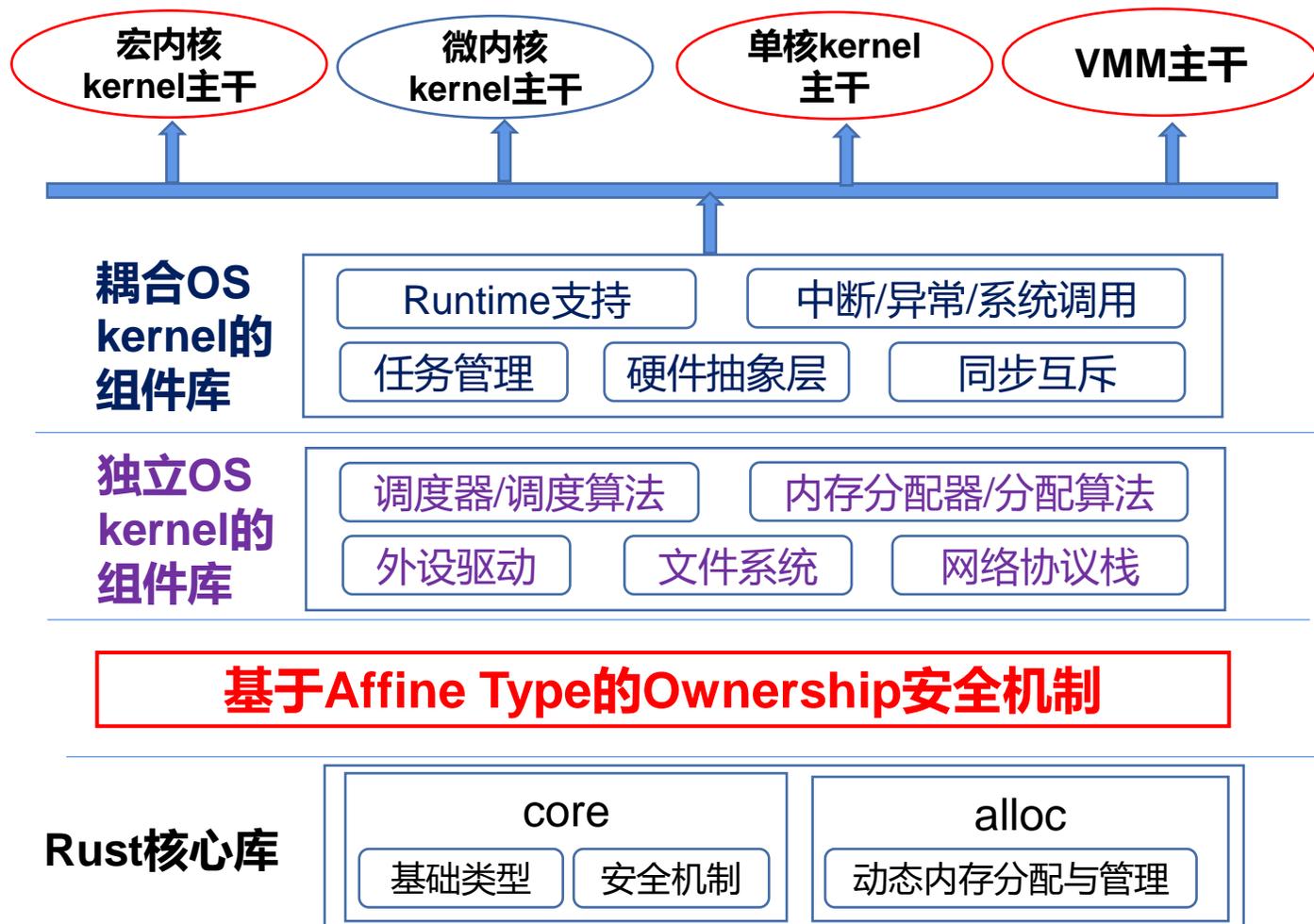
Substructural type system

主要研究内容及技术架构

操作系统架构设计:

细化内核模块属性, 形成单向依赖, 独立存在的内核模块

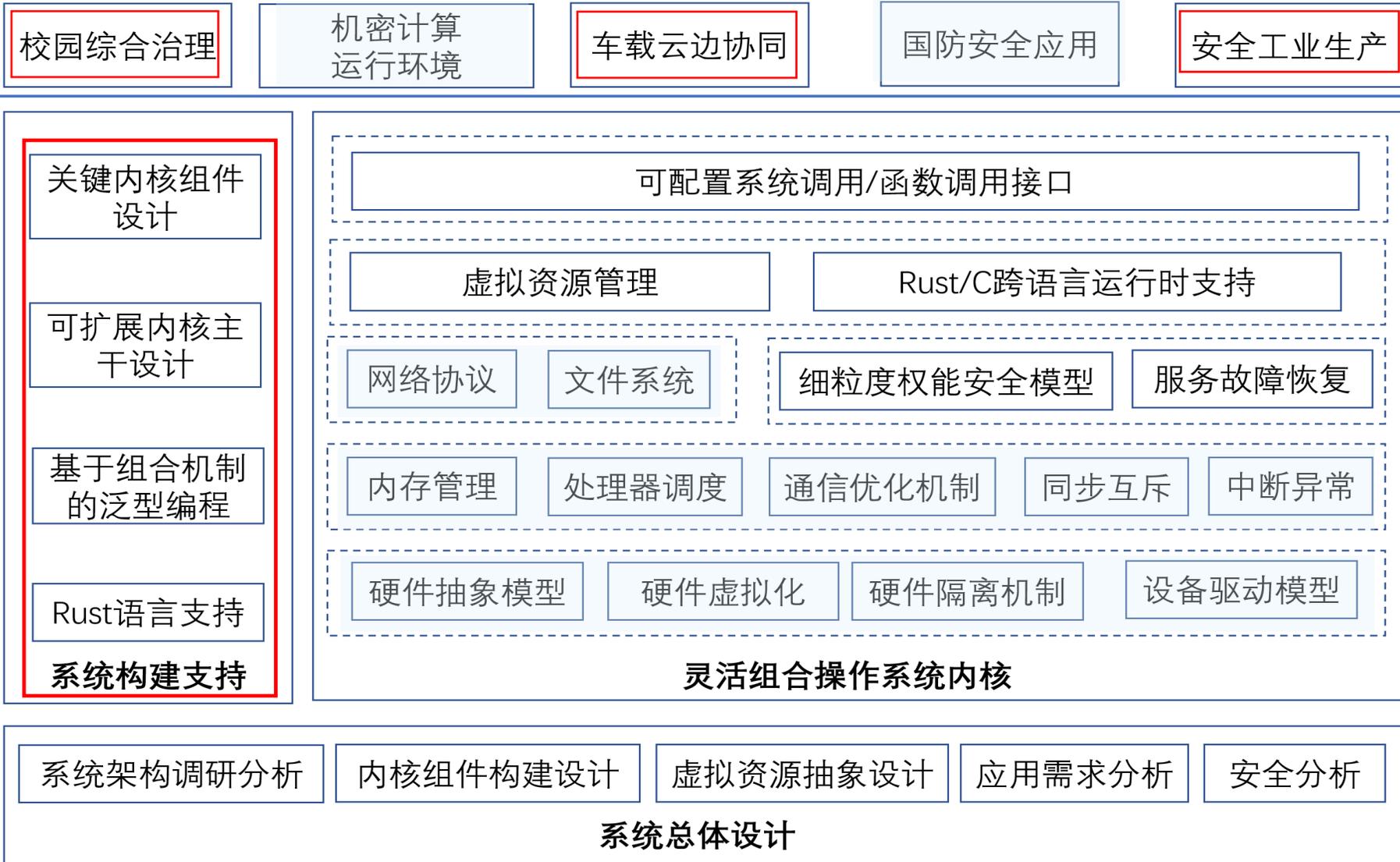
- 与OS无关的语言级核心库
- 独立OS的组件库
- 耦合OS的组件库
- 挂接组件的内核主干
- 独立OS的组件库无紧耦合依赖关系
- 基于配置和静态分析进行架构优化



主要研究内容及技术架构

■ 新型架构操作系统内核实现

分别从关键系统组件、操作系统内核主干、内核组合机制等多个角度分别进行研究

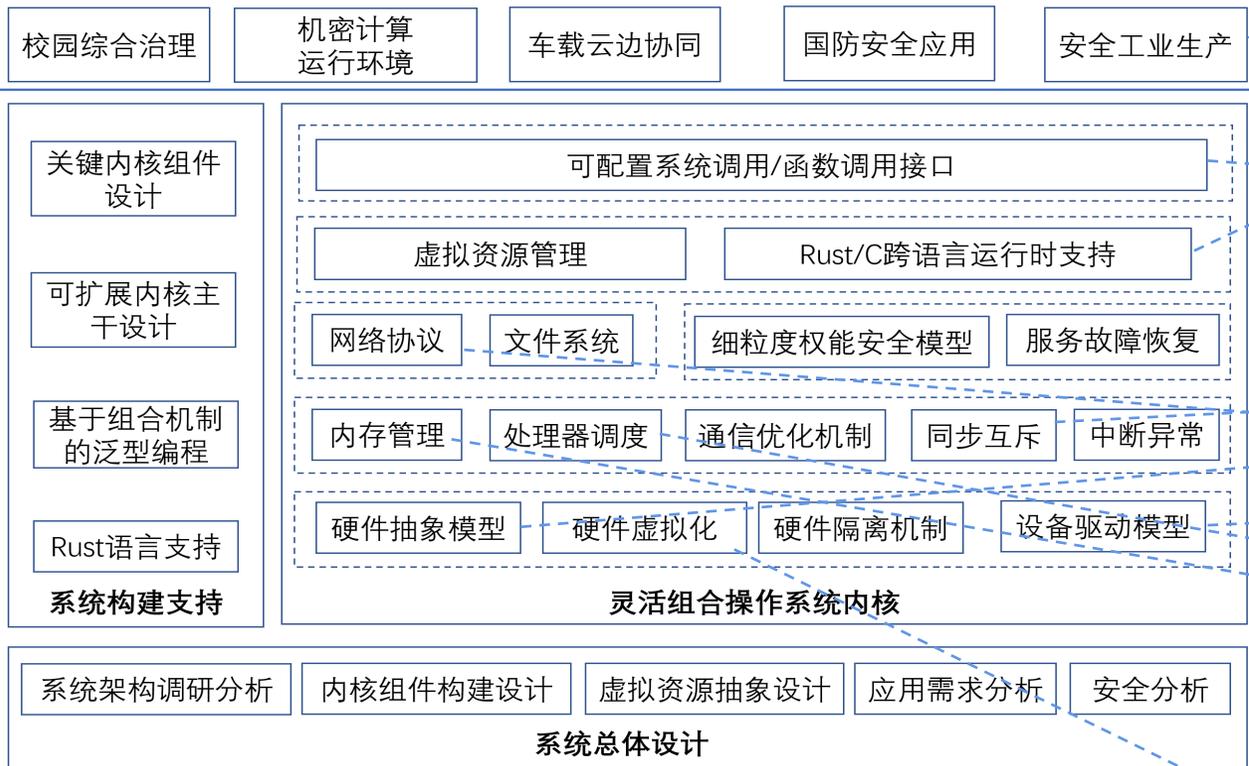


x86 ARM RISC-V LoongArch

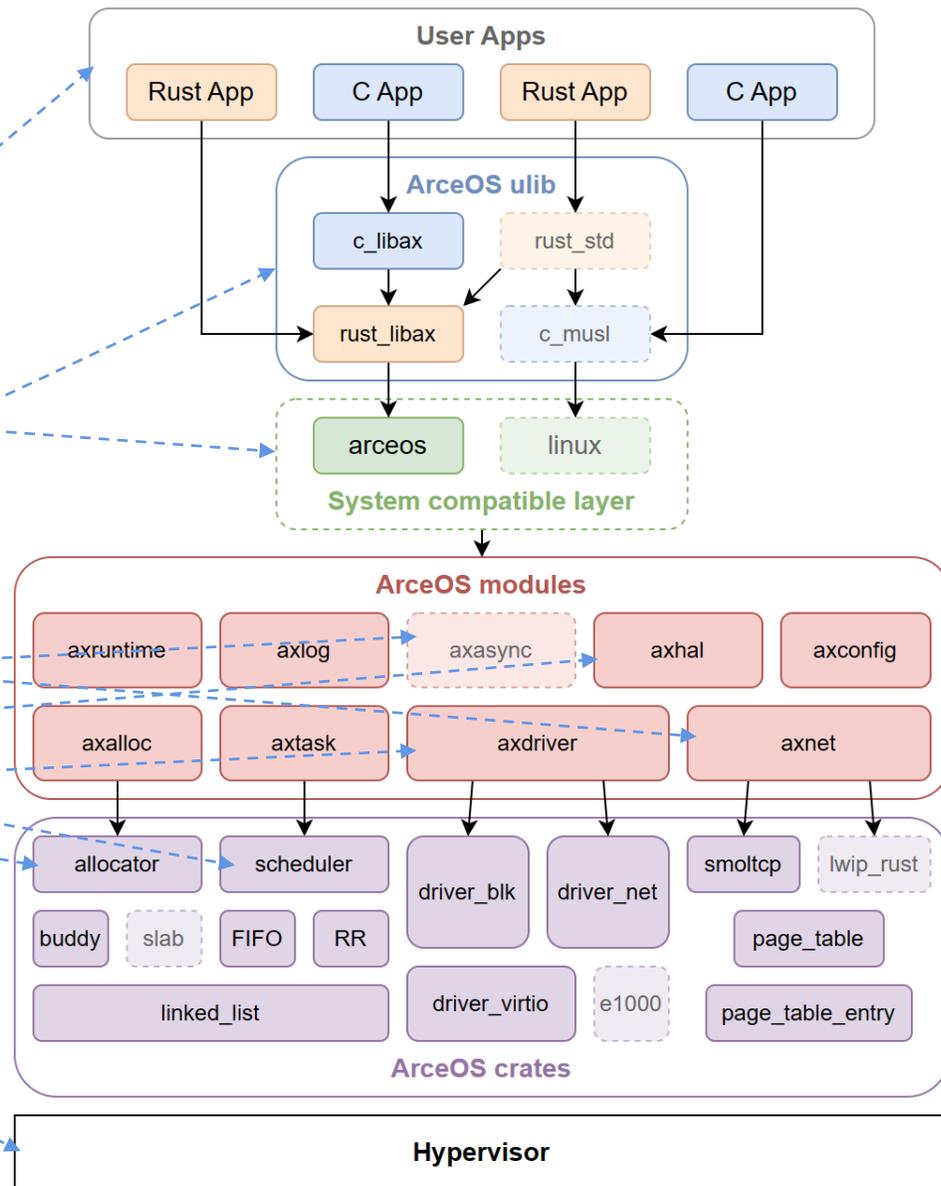
GPU NPU TPU

Devices

主要研究内容及技术架构



Unikernel架构模型



Unikernel架构的ArceOS内核

关键技术及指标

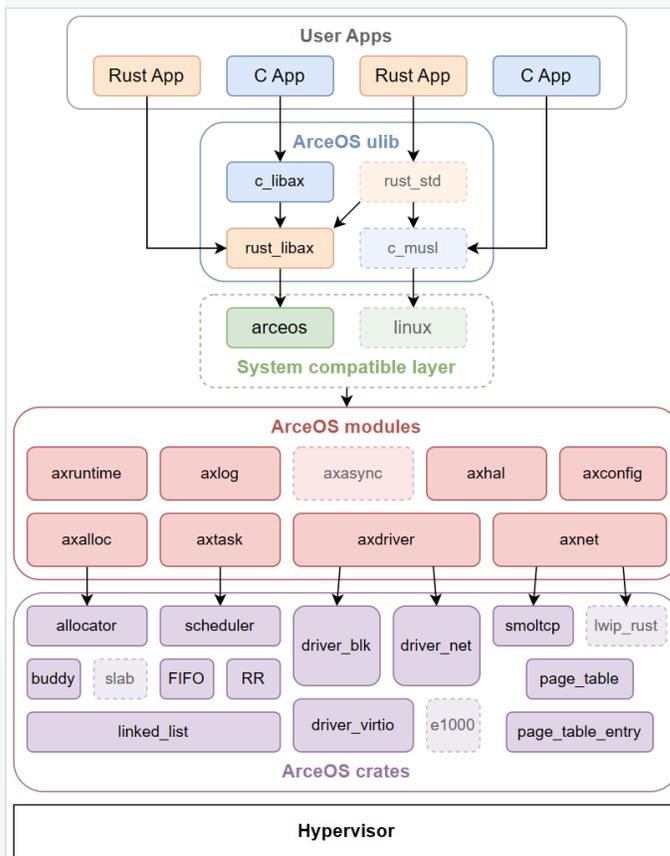
操作系统内核安全组合机制

支持最小权限的权能安全技术

- 基于组合模式和泛型模式进行功能扩展
- 特别是采用Rust语言的trait机制
- 实现操作系统内核组件在行为与数据上的分离
- 内核组件具有泛型编程的特征
- 通过Ownership机制的组合约束确保安全性

预期成果:

支持Hypervisor和kernel 形态, 支持微内核、宏内核、Unikernel等操作系统架构, 支持100个以上的Linux Syscall



ArceOS Modules

- **axalloc**: ArceOS global memory allocator.
- **axconfig**: Platform-specific constants and parameters.
- **axdisplay**: ArceOS graphics module.
- **axdriver**: ArceOS device drivers.
- **axfs**: ArceOS filesystem module.
- **axhal**: ArceOS hardware abstraction layer, prc.
- **axlog**: Macros for multi-level formatted logging.
- **axnet**: ArceOS network module.
- **axruntime**: Runtime library of ArceOS.
- **axsync**: ArceOS synchronization primitives.
- **axtask**: ArceOS task management module.

Crates

- **allocator**: Various allocator algorithms in a unified interface.
- **arm_gic**: ARM Generic Interrupt Controller (GIC) register definitions and base address.
- **axerrno**: Error code definition used by ArceOS.
- **axfs_devfs**: Device filesystem used by ArceOS.
- **axfs_vfs**: Virtual filesystem interfaces used by ArceOS.
- **axio**: `std::io`-like I/O traits for `no_std` environment.
- **capability**: Provide basic capability-based security.
- **crate_interface**: Provides a way to define an interface (trait) in a crate, but not in a module.
- **driver_block**: Common traits and types for block storage drivers.
- **driver_common**: Device driver interfaces used by ArceOS.
- **driver_display**: Common traits and types for graphics device drivers.
- **driver_net**: Common traits and types for network device (NIC) drivers.
- **driver_pci**: Structures and functions for PCI bus operations.
- **driver_virtio**: Wrappers of some devices in the `virtio-drivers` crate, that implement the virtio interface.
- **handler_table**: A lock-free table of event handlers. [crates.io v0.1.0](https://crates.io/v0.1.0)
- **kernel_guard**: RAIL wrappers to create a critical section with local IRQs or prc.
- **lazy_init**: A wrapper for lazy initialized values without concurrency safety but with lazy initialization.
- **linked_list**: Linked lists that supports arbitrary removal in constant time.
- **memory_addr**: Wrappers and helper functions for physical and virtual addresses.
- **page_table**: Generic page table structures for various hardware architectures.
- **page_table_entry**: Page table entry definition for various hardware architectures.
- **percpu**: Define and access per-CPU data structures.
- **percpu_macros**: Macros to define and access a per-CPU data structure.
- **ratio**: The type of ratios and related operations.
- **scheduler**: Various scheduler algorithms in a unified interface.
- **slab_allocator**: Slab allocator for `no_std` systems. Uses multiple slabs with 4096 bytes.
- **spinlock**: `no_std` spin lock implementation that can disable kernel local IRQs.
- **timer_list**: A list of timed events that will be triggered sequentially when the timer expires.
- **tuple_for_each**: Provides macros and methods to iterate over the fields of a tuple.

关键技术及指标

多态硬件抽象管理 驱动开发模型

- 硬件抽象层设计
- 硬件隔离支持
- 设备驱动模型
- 硬件抽象管理和开发模型
- 驱动程序的可靠性
- Aarch64/x86-64/RISCV-64/LoongArch64

预期成果:

支持国产自主可控的CPU/RISC-V/ARM, 支持网络/存储/AI加速芯片多种外设。

VirtIO-drivers-rs

crates.io v0.4.0 docs passing CI passing

VirtIO guest drivers in Rust. For `no_std` environment.

Device	Supported
Block	✓
Net	✓
GPU	✓
Input	✓
Console	✓
Socket	✓

Used by 332



Contributors 12



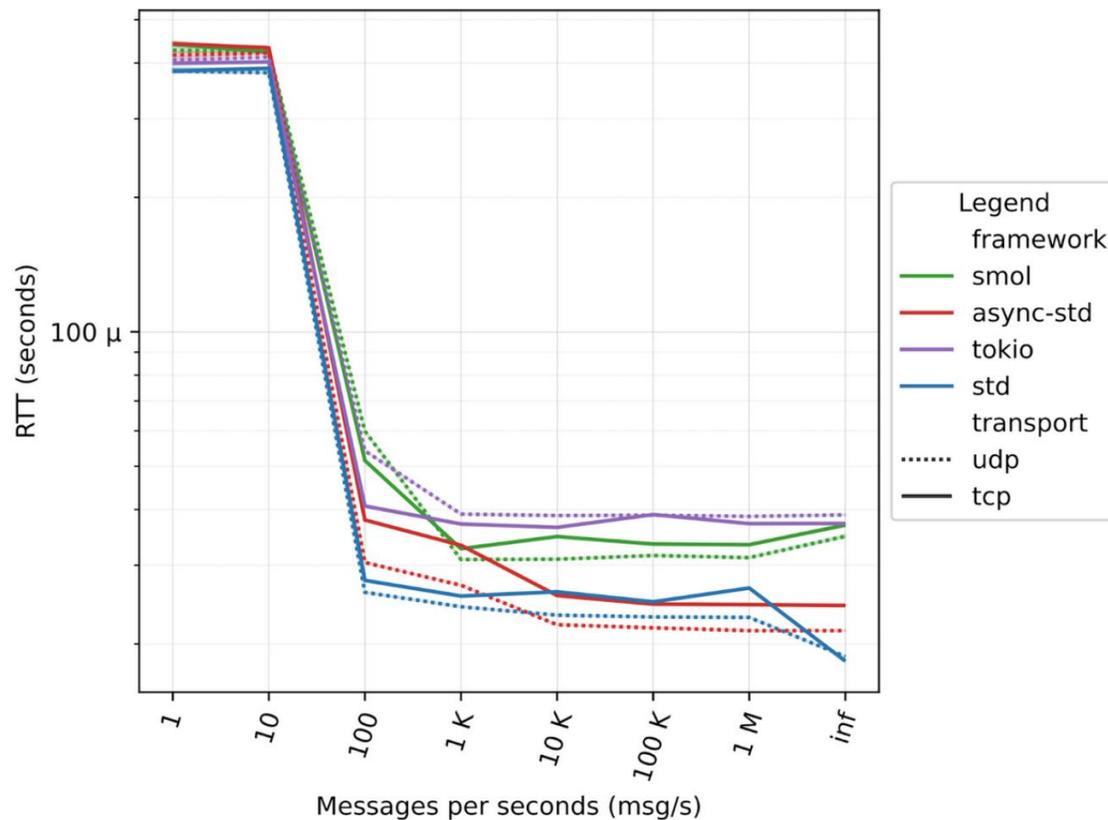
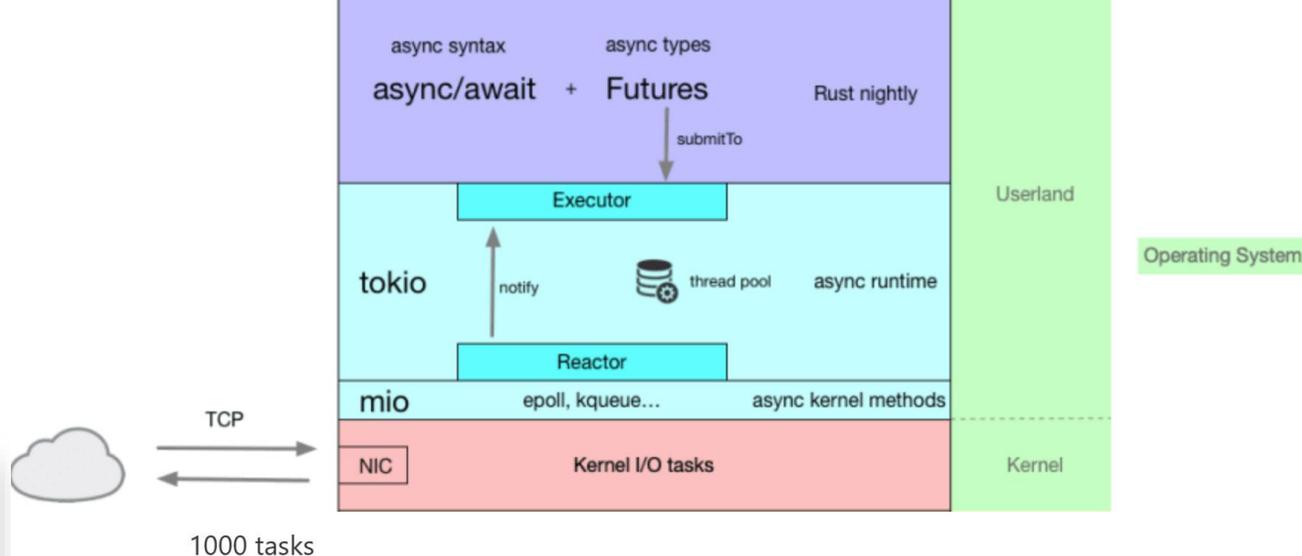
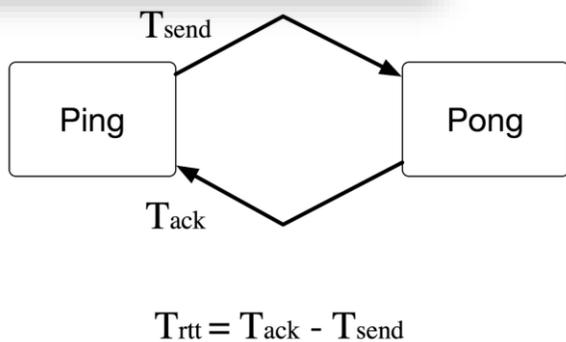
关键技术及指标

协程级异步调度机制

- 研究进程/线程/协程的嵌套调度机制
- 实现协程与线程之间的互协作机制
- 在内核协程设计方面，采用中断线程化设计
- 支持内核控制流在协程与线程之间进行切换

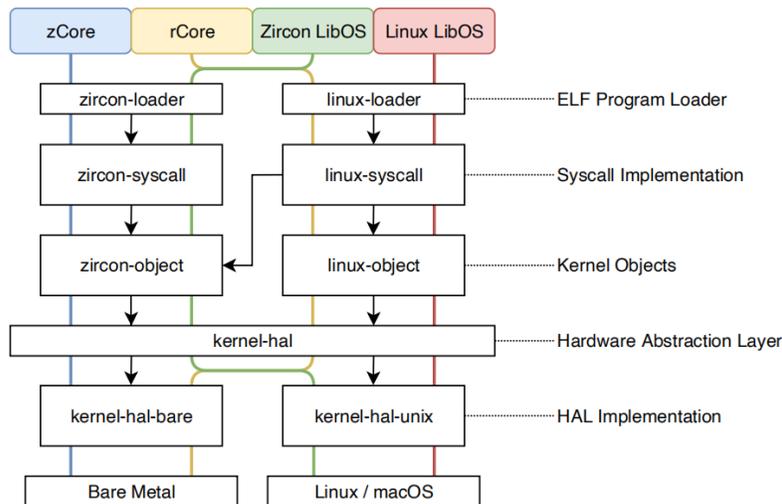
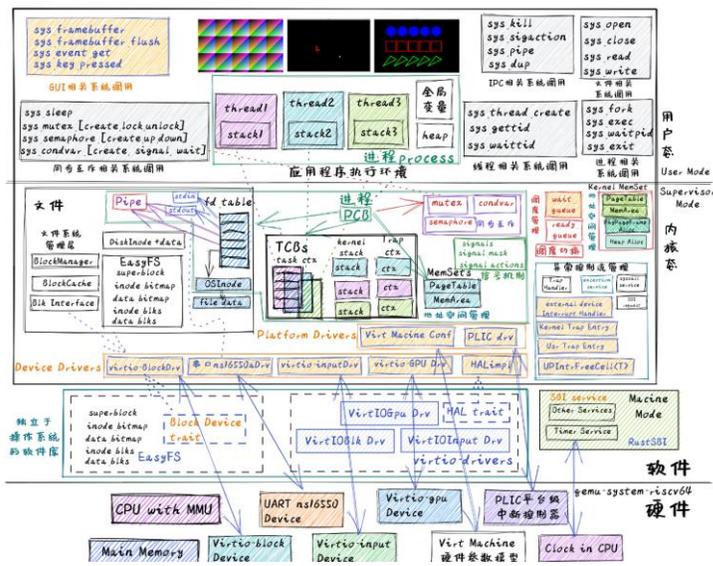
预期成果:

对redis/sqlite等典型Linux应用，比Linux Kernel在总体性能上超过50%。

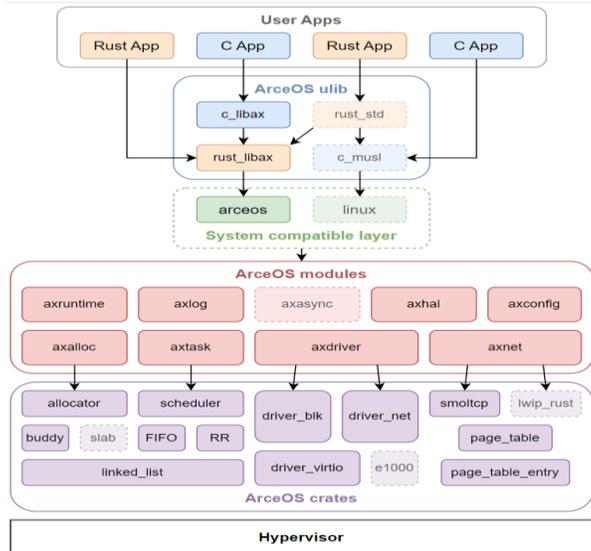


工作基础

uCore-tutorial → rCore-tutorial → rCore → zCore →



unikernel → multikernel



uCore-tutorial
教学操作系统

rCore-tutorial
国内第一个教学
Rust-based OS

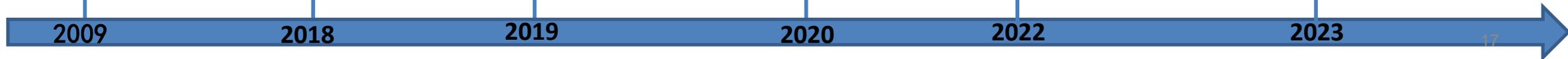
rCore
国内第一个科研
Rust-based OS

Occlum LibOS
与蚂蚁金服合作, 面向机密计算OS
zCore
多模态的Rust-based OS

HyperEnclave
与蚂蚁金服合作, 机密计算VMM
RVM
面向实时系统的Hypervisor

unikernel
泛在kernel
arceos
面向组合模式的OS框架

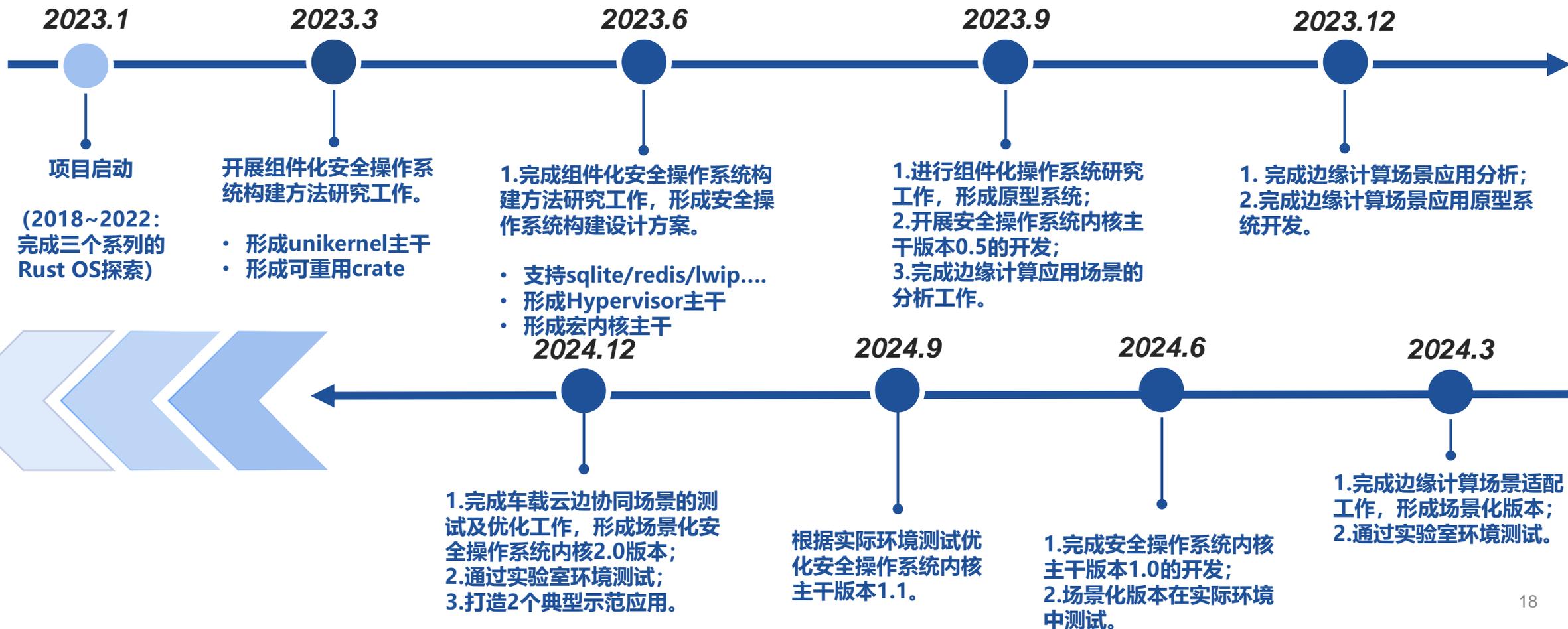
RVM-1.5
Hypervisor with linux



科研计划

清华大学 北京大学 中科院计算所 国汽智联创新中心 泉城实验室等

■ 以安全操作系统内核的研究和实现为主线，以应用场景的验证为辅助组织科研工作。



目录

- 一 | 整体规划
- 二 | 科研计划
- 三 | 落地规划

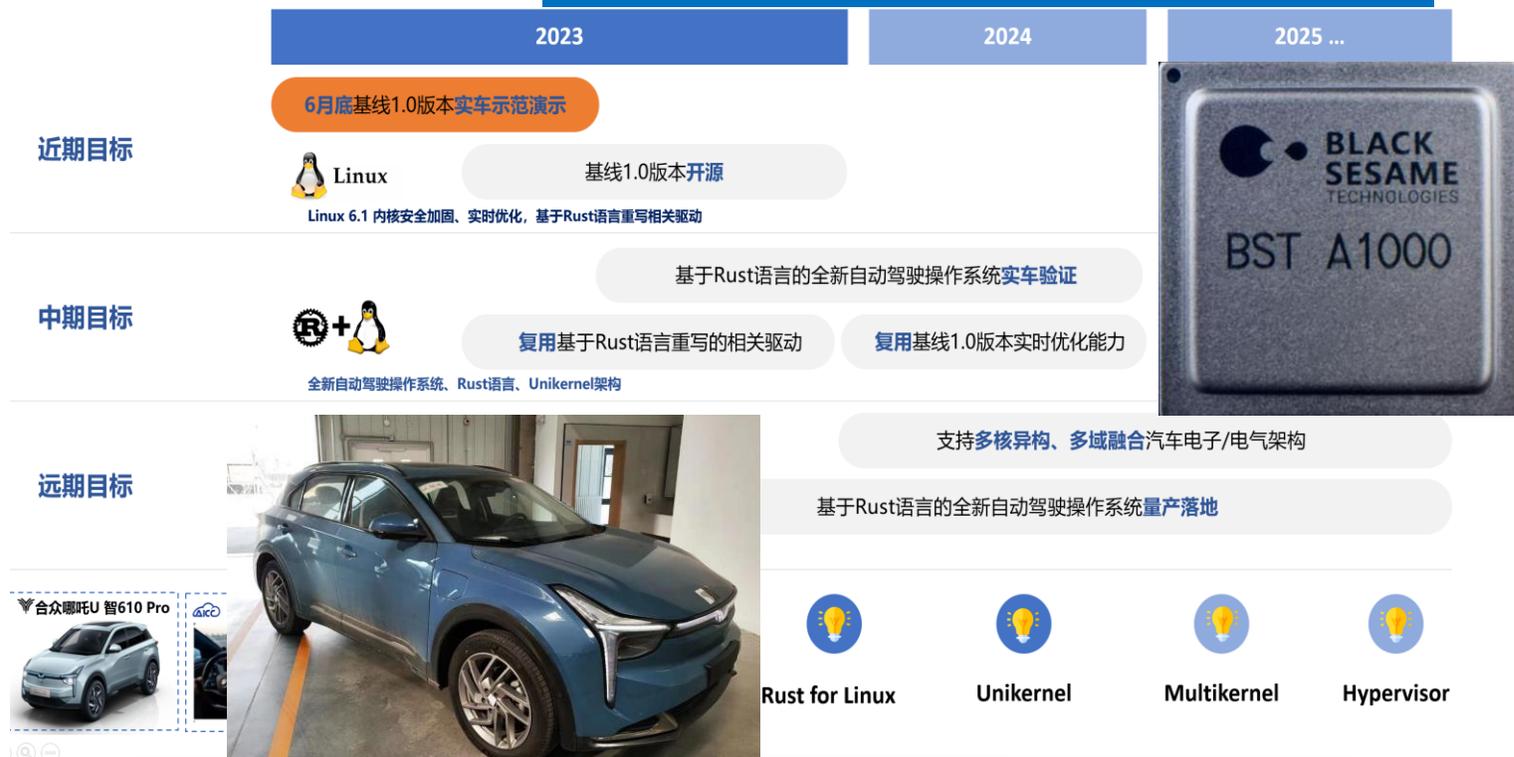


重点应用场景：智能网联车用操作系统内核

场景描述

基于**安全操作系统内核**，支持一端连接智能网联汽车整车企业，确保整车应用层面自动驾驶、智能网联、信息娱乐等功能的实现；一端连接芯片企业，确保不同功能的硬件资源可以充分适配调用。

CSAE-CCF-CICV操作系统联合实验室



- 1. 驱动代码移植&验证完成**
 - 从 Linux v5.10 到 v6.1
 - 总线及设备驱动40余种
- 2. 实车验证进行中**
 - 合众哪吒自动驾驶功能适配
- 3. 多单位联合开发**
 - 清华大学、电子科大
 - 黑芝麻、诚迈科技
 - 创新中心
- 4. 团队协同工作**
 - 建设代码GIT仓库
 - 编写开发文档
 - 技术经验分享



天目全数字实时仿真软件 SkyEye

SkyEye起源

源自于清华大学2003年陈渝老师发起的开源软件，被国内外个人爱好者和高校大量使用，且有相关操作系统书籍以及高校课程。

目前由康烁老师&迪捷软件负责商业化开发，广泛用于飞机、导弹、卫星和汽车等安全关键领域。

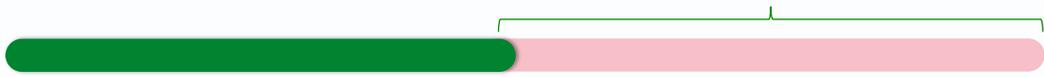
SkyEye介绍

- 基于可视化建模的硬件行为级仿真平台，支持用户通过拖拽的方式对硬件进行行为级别的仿真和建模。
- 仿真ARM、PowerPC、X86、DSP、Sparc、MIPS、龙芯、飞腾、国微等主流嵌入式硬件平台。
- 适配国内自主研发的操作系统天脉、锐华、翼辉、RT-Thread等。
- 通过利用基于LLVM的动态二进制翻译技术，使虚拟处理器运行速度可以达到2000MIPS以上。

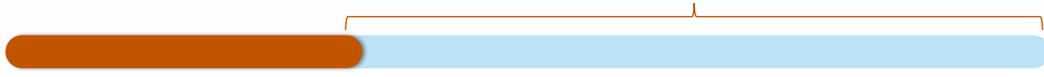
调试时间节省35%



测试成本减少50%



产品上市时间缩短66%

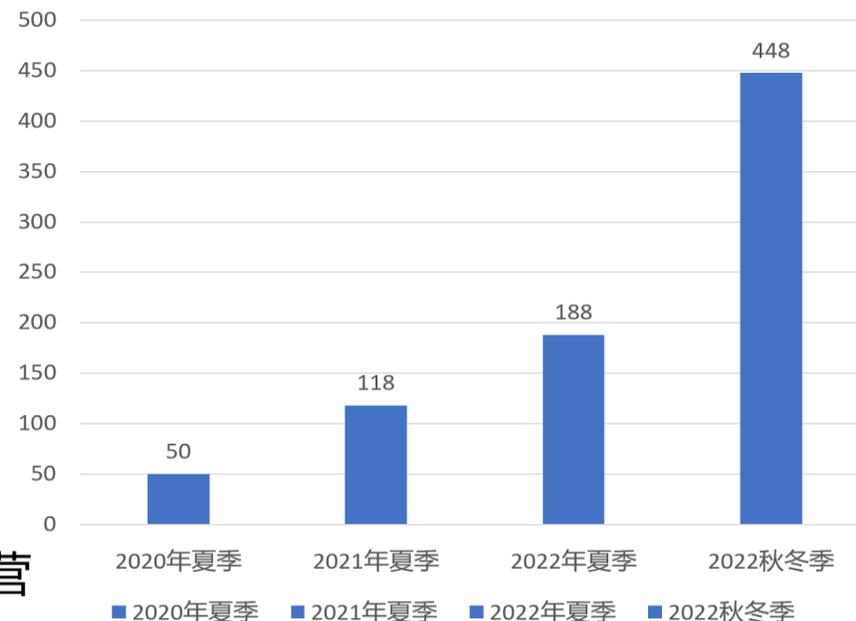


人才培养与团队建设



系统软件比赛 & 开源操作系统训练营

- 2020年
 - 2020.7.1 ~ 8.31 夏季训练营
- 2021年  CSCC 全国大学生计算机系统能力大赛
 - 2021.7.1 ~ 7.31 夏季训练营
- 2022年  第五届中国软件开源创新大赛 CSCC 全国大学生计算机系统能力大赛
 - 2022.7.3 ~ 8.31 夏季训练营
 - 2022.10.16 ~ 2023.2.1 秋冬季训练营



rCoreOS开源社区

探索基于Rust语言开发各种架构/形态的操作系统内核框架和可重用的内核组件，帮助应用和内核开发者快速搭建满足特定应用需求的高性能高安全操作系统内核。

- 想法来源: 基于泛型独立组件构建各种领域OS
- ArceOS: 组件化操作系统的初步探索

开源操作系统年度技术会议 OS2ATC

中国Linux内核开发者大会

人才培养与团队建设

- 开源OS训练营: <https://github.com/LearningOS>
- 开源OS科研项目: <https://github.com/rcore-os>
- 清华OS课件: <https://learningos.github.io/os-lectures>
- 全国大学生OS比赛: <https://os.educg.net>

产	国汽智联, 合众汽车, 黑芝麻智能, 蚂蚁, 地平线等
学	清华大学, 北京大学, 中科院计算所, 北航等
研	鹏城实验室, 启元实验室, 中关村实验室, 泉城实验室, 车用泛在操作系统联合实验室等
用	智能网联汽车操作系统, 工业安全操作系统, 机器人操作系统, 物联网AIoT操作系统等

- 自动驾驶OS开发训练营: <https://github.com/cicvedu>

The image shows a screenshot of the GitHub repository for 'The Learning&Training Hub of OS Kernel' and a diagram of the 'CICV OS 技术路线' (CICV OS Technology Roadmap).

GitHub Repository: The Learning&Training Hub of OS Kernel. It includes a README.md file with the following content:

the Learning&Training Hub of OS Kernel

2023开源操作系统训练营

- 2023.05.27: 在5月底前完成第二阶段 rCore OS实验的全部内容, 排行榜分数达到500分(满分)或类似的能力证明的同学, 可联系李明老师了解详情, 参加内容丰富的线下实习计划。实习地点在北京/济南。
- 2023.05.07: 部分已经完成训练营第二阶段训练的同学, 可联系李明老师, 与全国的学生/工程师一起参加各种有趣或挑战性的小项目
- 2023.04.18: 大家在完成训练营第二阶段的训练后, 如果还有兴趣探索基于Rust的操作系统, 欢迎参与这些实验小项目。
- 2023.04.01: 2023年春夏季开源操作系统训练营正式启动。
- 2023.03.25: 年春夏季开源操作系统训练营启动&报名交流会(腾讯会议号: 231-190-126)
- 2023.03.23: 请参加的同学填写训练营报名登记表, 并加入相关微信群。目前群人数已经200人, 入群请联系李明老师(微信id: limingth)
- 2023.03.19: 2023年春夏季开源操作系统训练营启动&报名交流会, 会议时间: 2023/03/25 11:00-12:00, #腾讯会议号: 231-190-126 会议密码: 0325, 请报名的同学参加
- 2023.02.17: 2023开源操作系统训练营准备中
- 2023.02.05: OS课程的部分大实验内容, 组件化操作系统设计与实现的相关题目, 欢迎感兴趣的朋友联系我们

Self Learning

CICV OS 技术路线

The diagram illustrates the technology roadmap for CICV OS, structured into layers:

- 应用软件 (Application Software):** Includes 自动驾驶 (Autonomous Driving), 工业控制 (Industrial Control), 机器人 (Robotics), 无人机 (Drones), 车联网 (V2X), and 控制执行 (Control/Execution).
- 功能软件 (Functional Software):** Includes 自动驾驶通用框架模块 (感知、规划、控制) (Autonomous Driving General Framework Modules: Perception, Planning, Control), 深度学习和视觉模块 (Deep Learning and Vision Modules), 传感器模块 (Sensor Modules), 网联模块 (V2X Modules), and 云控模块 (Cloud Control Modules).
- 系统软件 (System Software):** Includes 管理平面和实时控制平面 (Management and Real-time Control Planes), 分布式通信 (Distributed Communication), AUTOSAR RTE, Linux等 (Linux, etc.), Linux优化/其他安全实时内核 (Linux Optimization/Other Safe Real-time Kernels), Safety RTOS, BSW, Hypervisor及BSP Drivers, and MCAL.
- 硬件平台 (Hardware Platform):** Includes AI单元-GPU/FPGA/ASIC AI芯片+CPU (AI Units-GPU/FPGA/ASIC AI Chips+CPU), 计算单元-多核多CPU (Computing Units-Multi-core Multi-CPU), and 控制单元-MCU (Control Units-MCU).

The diagram also highlights a **安全体系 (Safety System)** on the right side, involving 工具链 (开发、仿真、调试、测试、验证) (Toolchain: Development, Simulation, Debugging, Testing, Verification) and 安全体系 (Safety System).

谢谢！
请专家批评指正