

字节大规模内核版本迁移实践

字节跳动
贺中坤



Why

内核版本在社区支持时间有限
硬件迭代
软件新feature带来优化

What

保证稳定性的无感升级

How

怎样才能做到

通算集群的挑战

- 业务场景复杂
- 混部后单机压力大，触发边界问题
- 内核版本RC后的测试难以覆盖全
- 迁移中途遇见问题回退版本

稳定性要求高，迭代慢

新增智算集群挑战

- 智算逐步超过通算
- 围绕者高速的GPU逐步重构硬件和基础架构
 - 更强大单节点
 - 更高速的互联
 - 硬件迭代速度快

对稳定性和迭代速度有更高的要求

目录

C O N T E N T S

01 内核版本升级挑战

02 稳定性收敛机制

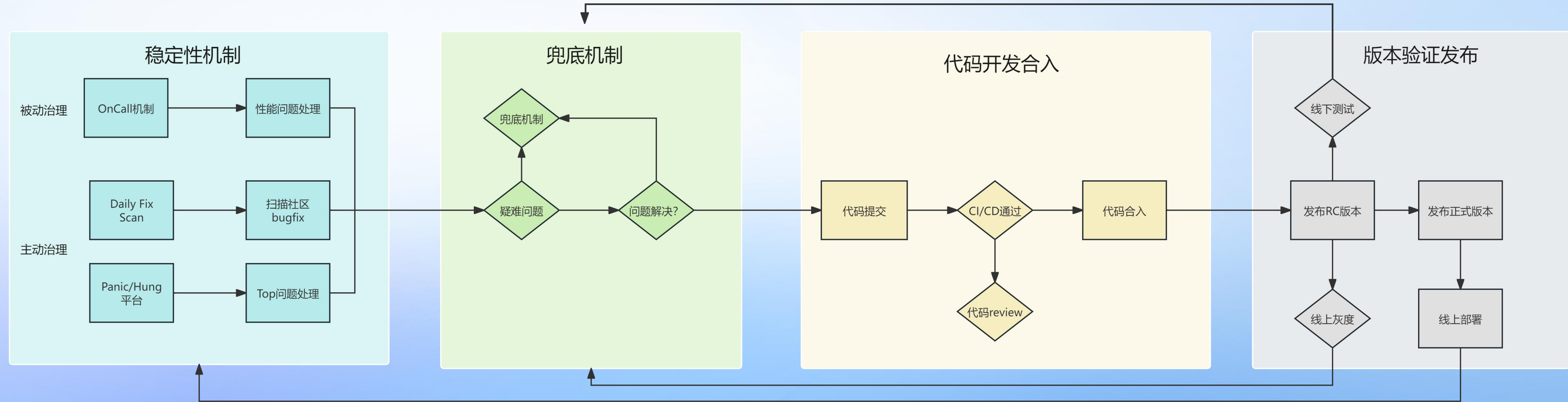
03 兜底能力建设

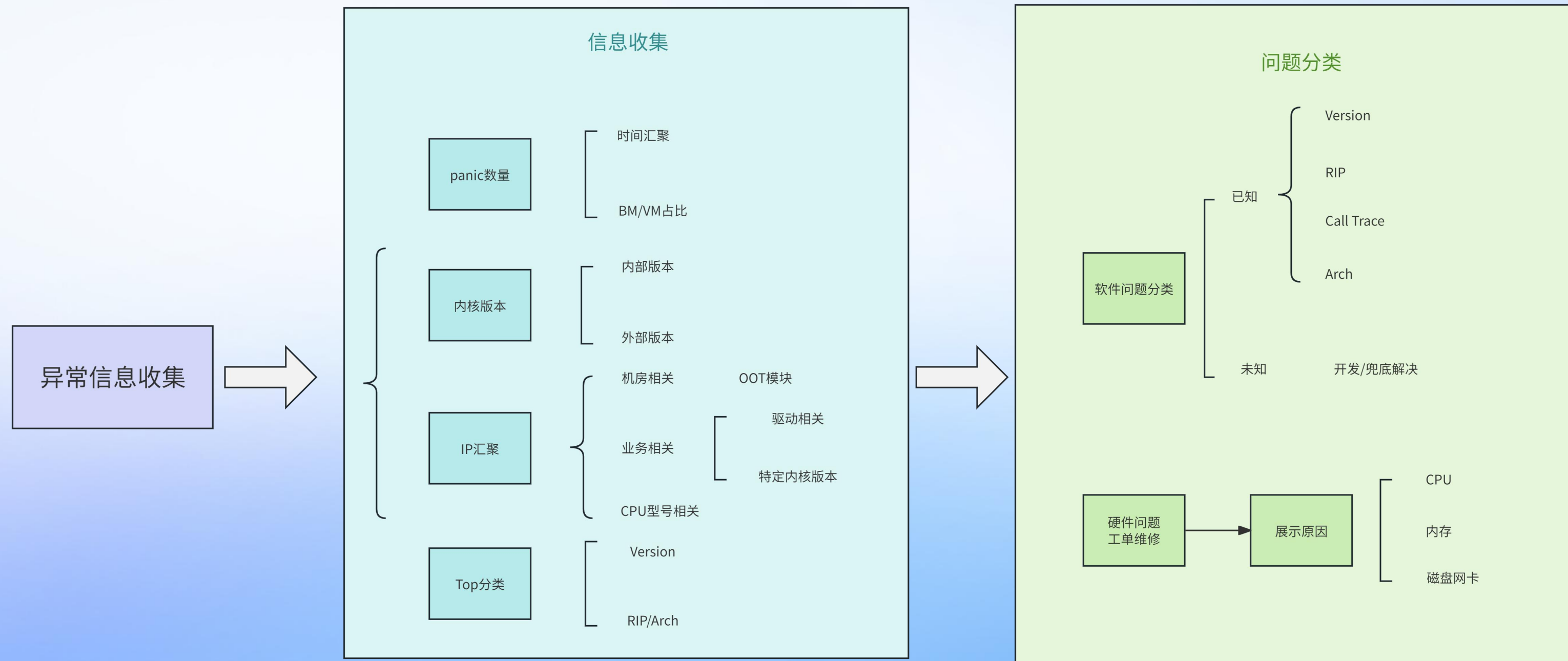
04 代码开发合入机制

05 发版验证机制

06 实践成果







稳定性建设-AI检索问题

DeepSeek

1

通过rip 解析函数名和内核版本

2

Dmesg 解析报错文件名 (by AI)

3

社区Master分支的 git commit log 查询近几年
该文件修改commit

4

逐个分析这些commit 是否解决该问题 (by AI)

```
[2205663.312099] RIP: 0010:ovl_dentry_upper+0x9/0x20 [overlay]
[2205663.317222] Code: 00 48 8b 53 30 5b 48 8b 92 48 02 00 00 48 89 55 08 5d c3 66 66 2e 0f 1f 84 00 00 00 00 00
[2205663.381971] Call Trace:
[2205663.383259] ovl_dentry_revalidate_common+0x19/0x80 [overlay]
[2205663.385014] lookup_fast+0x103/0x130
[2205663.386728] walk_component+0x40/0x1a0
[2205663.388082] link_path_walk.part.45+0x226/0x340
[2205663.389660] ? path_init+0x2ef/0x360
[2205663.391356] path_openat+0xb1/0x1070
[2205663.396265] do_filp_open+0x93/0x100
```

深度检索社区相关patch

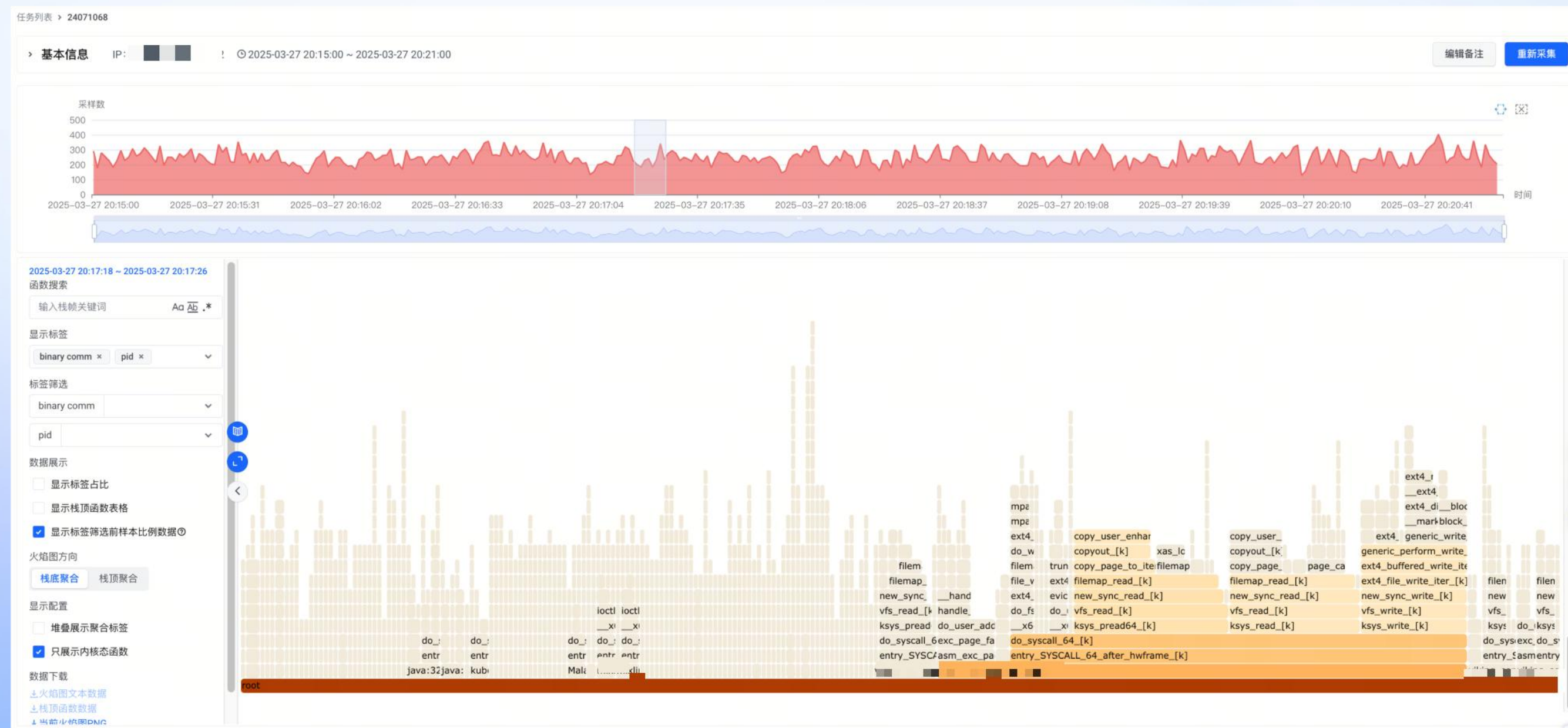
1 条回复

回复 贺中坤: [2205663.312099] RIP: 0010:ovl_dentry_upper+0x9/0x20 [overlay] [2205663.317222] Cod...

本次检索召回7个相关邮件组, 分别是:

- 1.ovl: Check for NULL d_inode() in ovl_dentry_upper()
该commit直接修复了ovl_dentry_upper()中d_inode()为NULL的情况, 与用户故障堆栈中的ovl_dentry_upper函数调用完全匹配。相关性得分9
- 2.ovl: fix null pointer dereference in ovl_get_acl_rcu()
该commit修复了RCU模式下获取ACL时的空指针问题, 与用户报告的ovl_dentry_upper故障有部分关联性, 可能间接解决类似问题。相关性得分8
- 3.ovl: fix null pointer dereference in ovl_permission()
该commit修复了权限检查时的空指针问题, 与用户报告的ovl_dentry_upper故障有部分关联性, 可能间接解决类似问题。相关性得分8
- 4.ovl: let helper ovl_i_path_real() return the realinode
该commit优化了获取真实inode的辅助函数, 与用户报告的ovl_dentry_upper故障有部分关联性, 可能间接解决类似问题。相关性得分8
- 5.acl: move idmapped mount fixup into vfs_{g,s}etxattr()
该commit修复了POSIX ACLs在idmapped mounts上的问题, 与故障现场中ovl_dentry_upper的调用栈相关, 可能解决用户问题。相关性得分8
- 6.ovl: fetch inode once in ovl_dentry_revalidate_common()
该commit优化了ovl_dentry_revalidate_common()中的inode获取, 与用户报告的故障堆栈中的ovl_dentry_revalidate_common函数直接相关, 可能解决相关问题。相关性得分8
- 7.ovl: fix use inode directly in rcu-walk mode

稳定性建设-debug



perf的缺点是什么？

1：无法常态化采集

2：数据量大

3：无法采集用户态栈

如何debug 偶现的性能抖动问题？

典型业务场景	机器规格	活跃CPU使用量 (关联CPU占用)	活跃elf数量 (关联内存占用)	CPU 占用	内存 占用	存储 占用
clickhouse独占机器	128C / 1.5T	67C	253	0.13C	2.5G	100M / 天
mysql独占机器	128C / 1.5T	67C	253	0.13C	2.5G	100M / 天
容器场景(单机30 Pod)	128C / 1.5T	67C	253	0.13C	2.5G	100M / 天
容器场景(单机1000 Pod)	128C / 1.5T	67C	253	0.13C	2.5G	100M / 天

持续的profiling

1：开销低

2：包含内核用户态

3：存储占用低

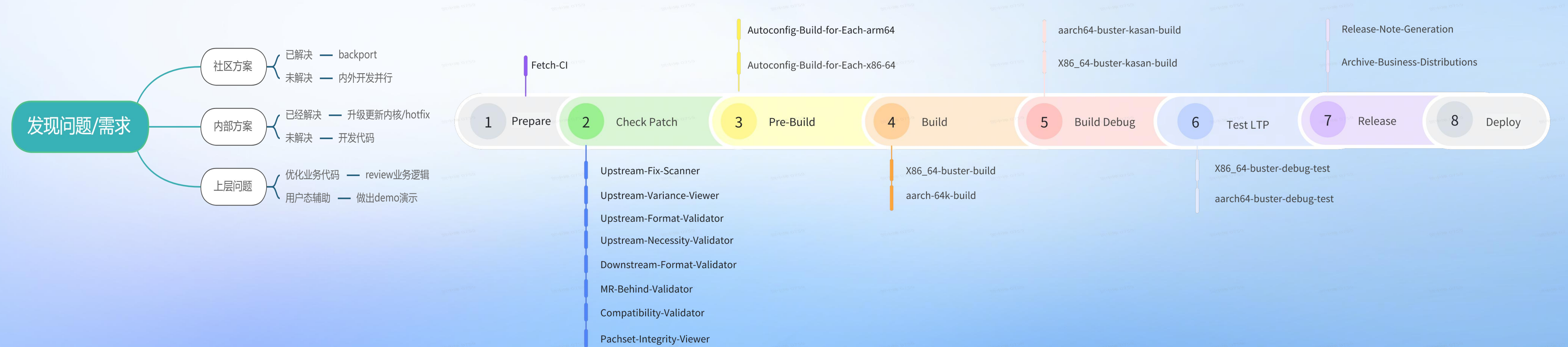
内核技术方向

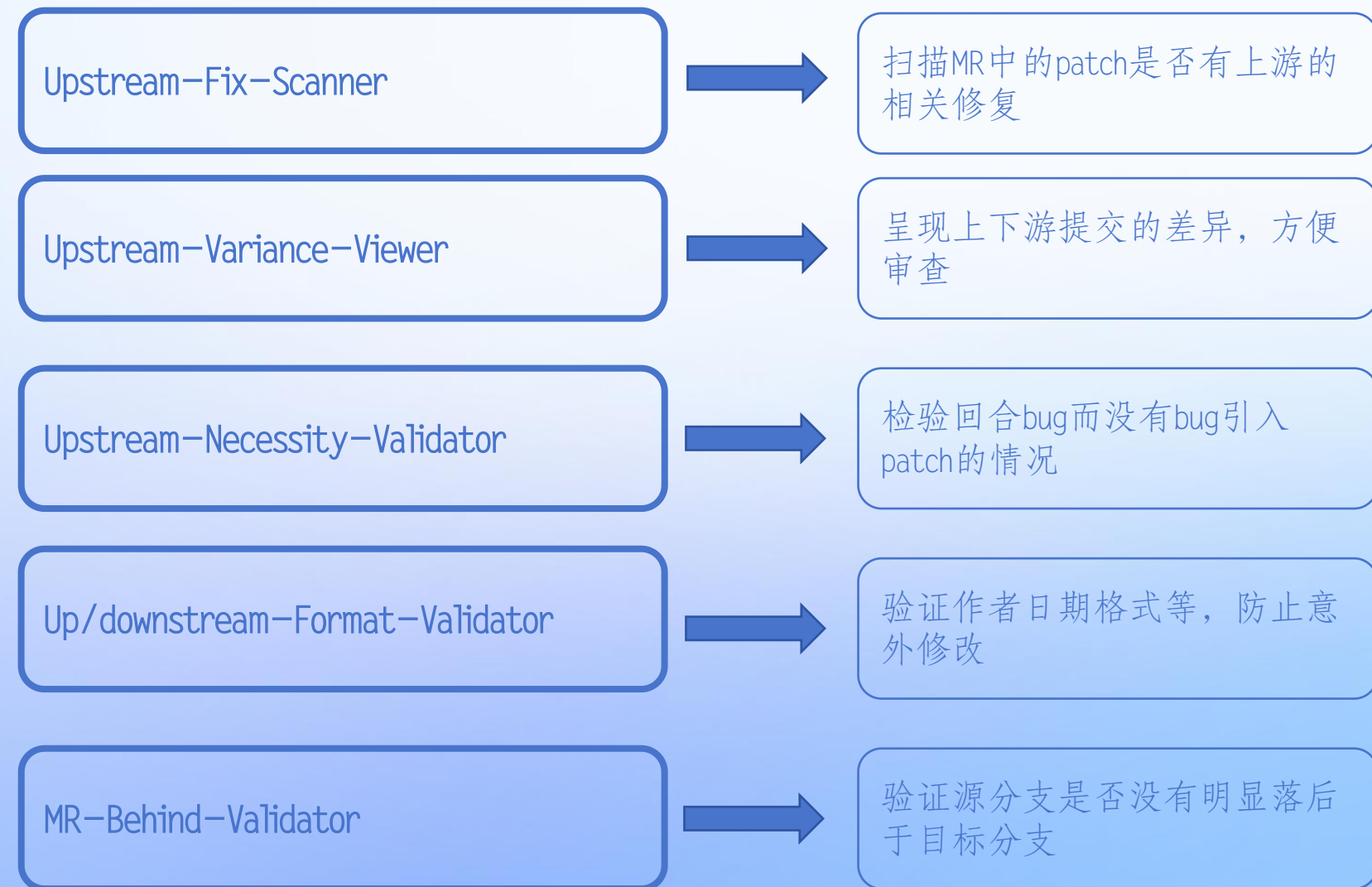
- 追随社区，回馈社区
- 方案推开源，尽量少做闭源代码
 - 长期技术债减少
- 构建调度/内存/存储/网络等方向专家团队

兜底还得依靠人

- ByteDance 在Linux社区贡献patches 960+
- 多个重要的feature合入Linux社区
- Maintainer/Reviewer 3位

代码开发合入机制



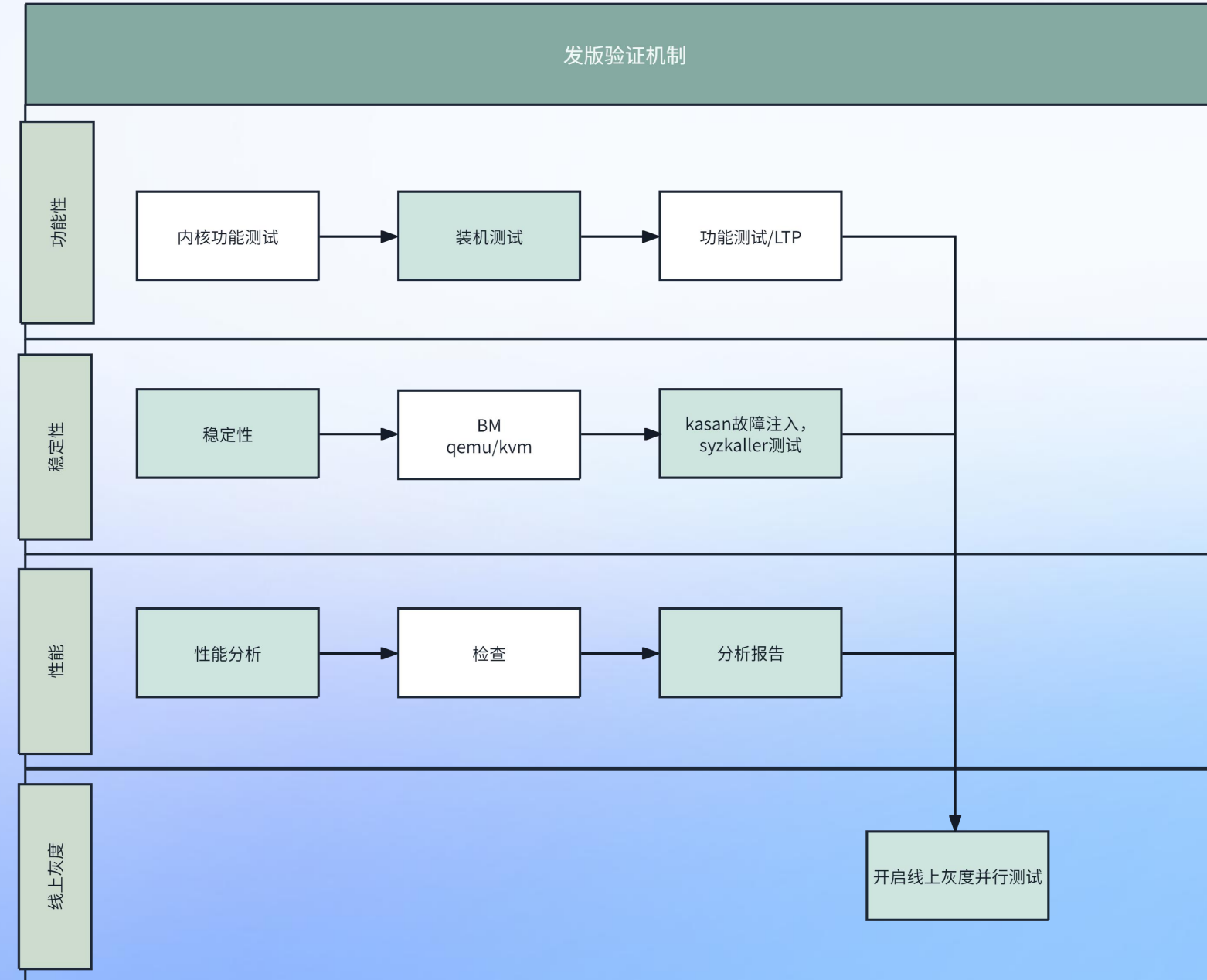


Daily Fix San

- 扫描fix标题的patch
- 使用AI处理不带fix标签的patch扫描，解决patch识别不准确的问题

识别出待合入的bufix，再相关领域进行筛选

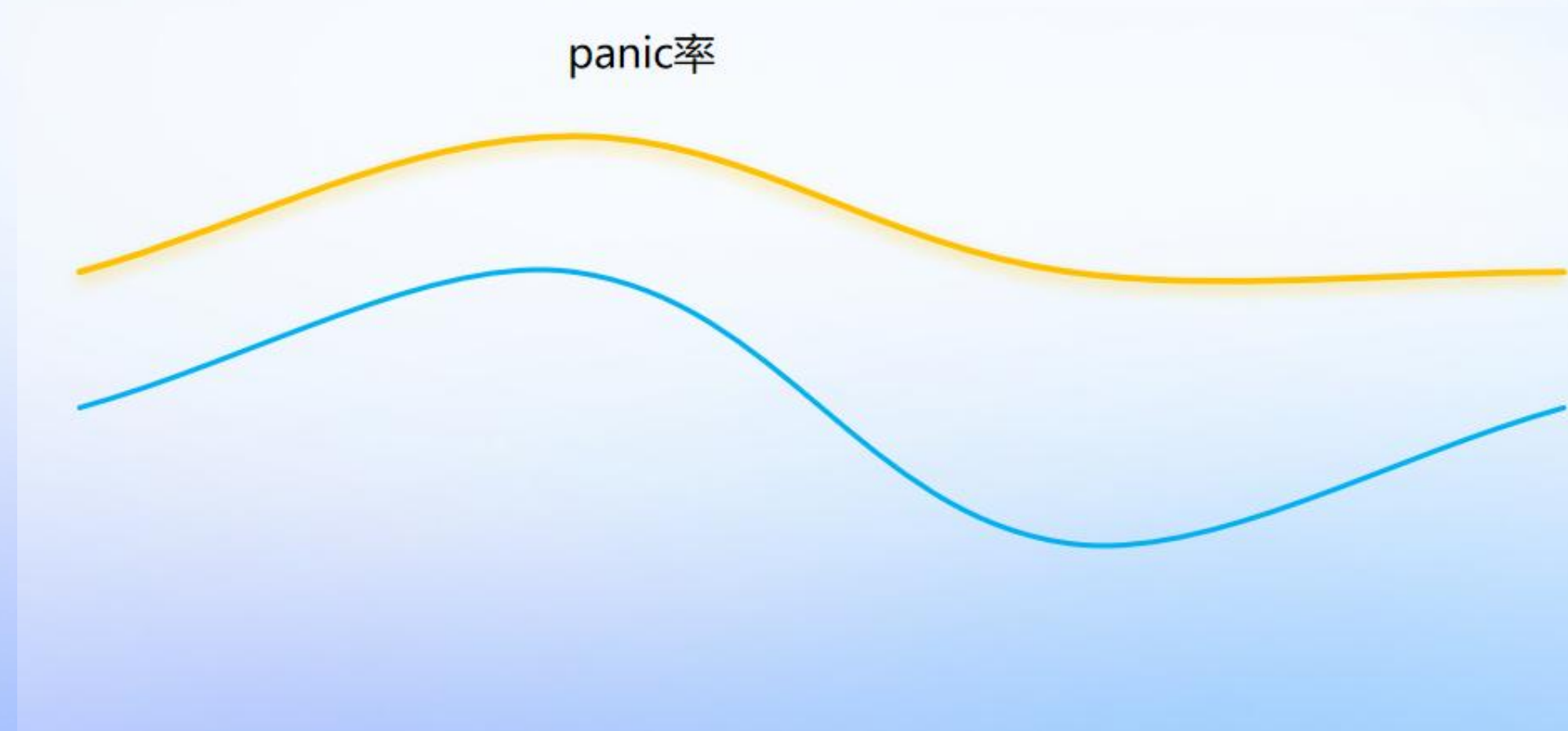
发版验证机制



6

实践成果展示





1: panic降低

2: hungtask降低

3: 问题定位效率大幅提升

问题现象

- 混部场景
- 低优业务获得cpu时间较少
- 持有全局锁(cgroup)
- 得不到运行后整机卡顿

解决方案

- 限时标记cfs_rq throttled但并不将它从rq上摘除
- task在返回用户态之前将自己从rq上摘除, 这时并不持有内核资源
- 方案已经合入6.18

问题现象

- 内存申请路径
- 持有全局资源(jbd2 handle)
- 全局或者cgroup内存不足
- loop in memory reclaiming

解决方案

- add task flag ACCOUNTFORCE
- cgroup内存不足时force charge
- 强制charge部分在退出到用户态之前回收, 不持有锁
- 准备推动方案force charge,回退到用户态时recalim和OOM
- 全局OOM场景没有较好的办法

问题现象

- 管控程序访问cpu.stat频率高
- cgroup_rstat_flush_locked
- 5.4 rstat延时 250us
- 5.15 rstat延时 7ms?
- cgroup v2关中断时间太长,影响网络收发包

解决方案

- 社区进行多年优化
 - global->sigle memcg flush
 - flush context sleep able
 - 全局rstat锁拆分到sub sys
- 5.15方案是 mutex global flush, 性能与准确度平衡

问题现象

- 用户态分配Jemalloc/TCMalloc
- 依赖 madvise系统调用释放物理内存
- 用户地址空间大规模膨胀
- PTE 页表占用约为虚拟地址空间大小的 2%。
- 线上发现数十GB的页表

解决方案

- page table相关的基础设施进行了进一步改造[1]
- 实现PTE页表的同步回收[2]
- 为修复其产生的UAF问题，再次对基础设施进行了改造[3]

[1] <https://lore.kernel.org/all/cover.1727332572.git.zhengqi.arch@bytedance.com/#r>

[2] <https://lore.kernel.org/all/cover.1733305182.git.zhengqi.arch@bytedance.com/>

[3] <https://lore.kernel.org/all/cover.1736317725.git.zhengqi.arch@bytedance.com/#r>



第二十届中国 Linux 内核开发者大会

📍 2025年11月1日 🕒 中国·深圳

THANKS





第二届中国Linux内核开发者大会

🕒 2025年11月1日 📍 中国·深圳

赞助单位



支持单位



支持媒体及社区

