

# New Features

字节跳动基础架构部内核与虚拟化团队  
祝嘉

## 实践与探索——New Features

01

**Failover**

02

EROFS Shared  
Domain

03

Daemonless

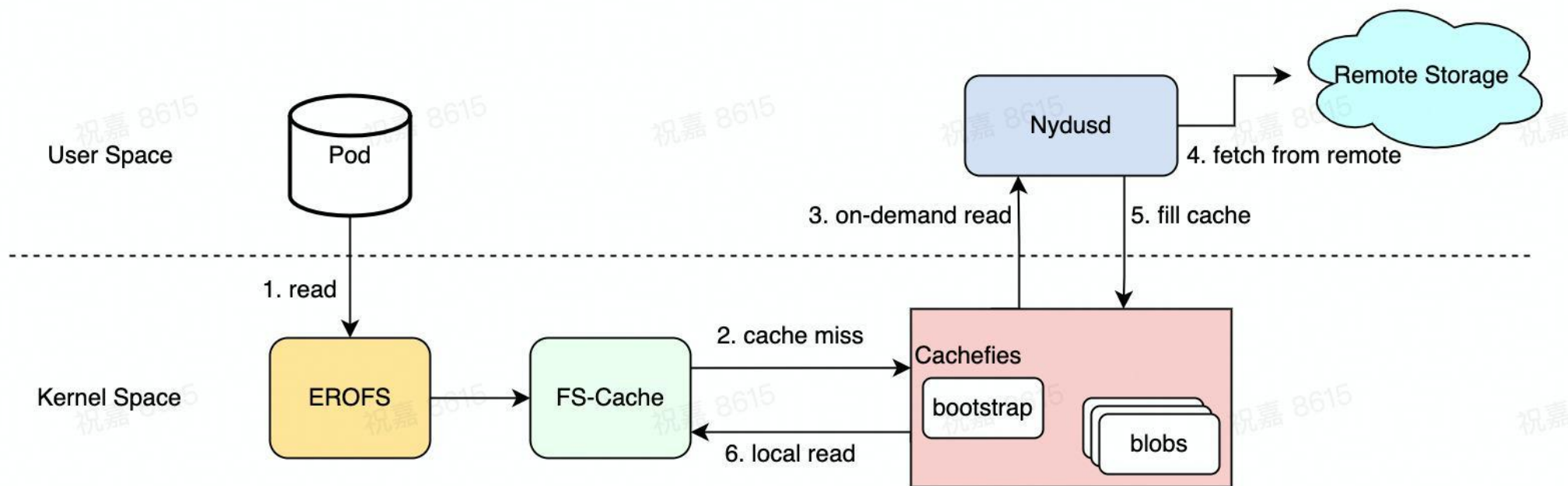
## 稳定性提升：User Daemon Failover

背景：

- 用户态进程crash后，IO请求将返回错误，影响容器中运行的业务

预期目标：容器不感知用户态进程crash/restart

- In-Flight IO不丢失
- User Daemon crash后 ~ 重启前 产生的IO不丢失
- IO在有限时间内恢复



# 稳定性提升：User Daemon Failover

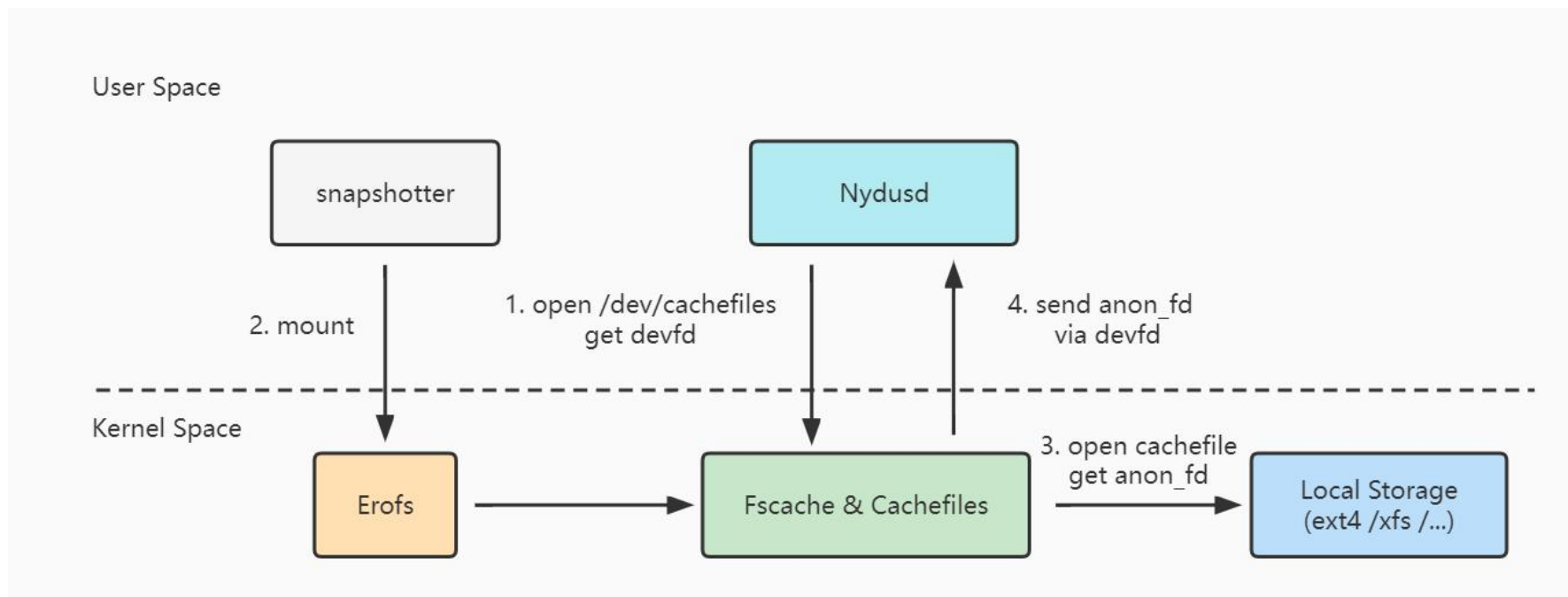
## 方案设计

- 异常监控
  - 用户态 supervisor 监控 user daemon 状态，触发故障恢复流程
- 资源保持与恢复
  - 控制面： /dev/cachefiles fd的保持和传递 (UDS)
  - 数据面： Reopen mechanism： 由按需IO驱动重新发送 anonymous fd
  - cachefiles 'restore' cmd： 触发恢复In-Flight IO

## 稳定性提升：User Daemon Failover

fd资源保持与恢复

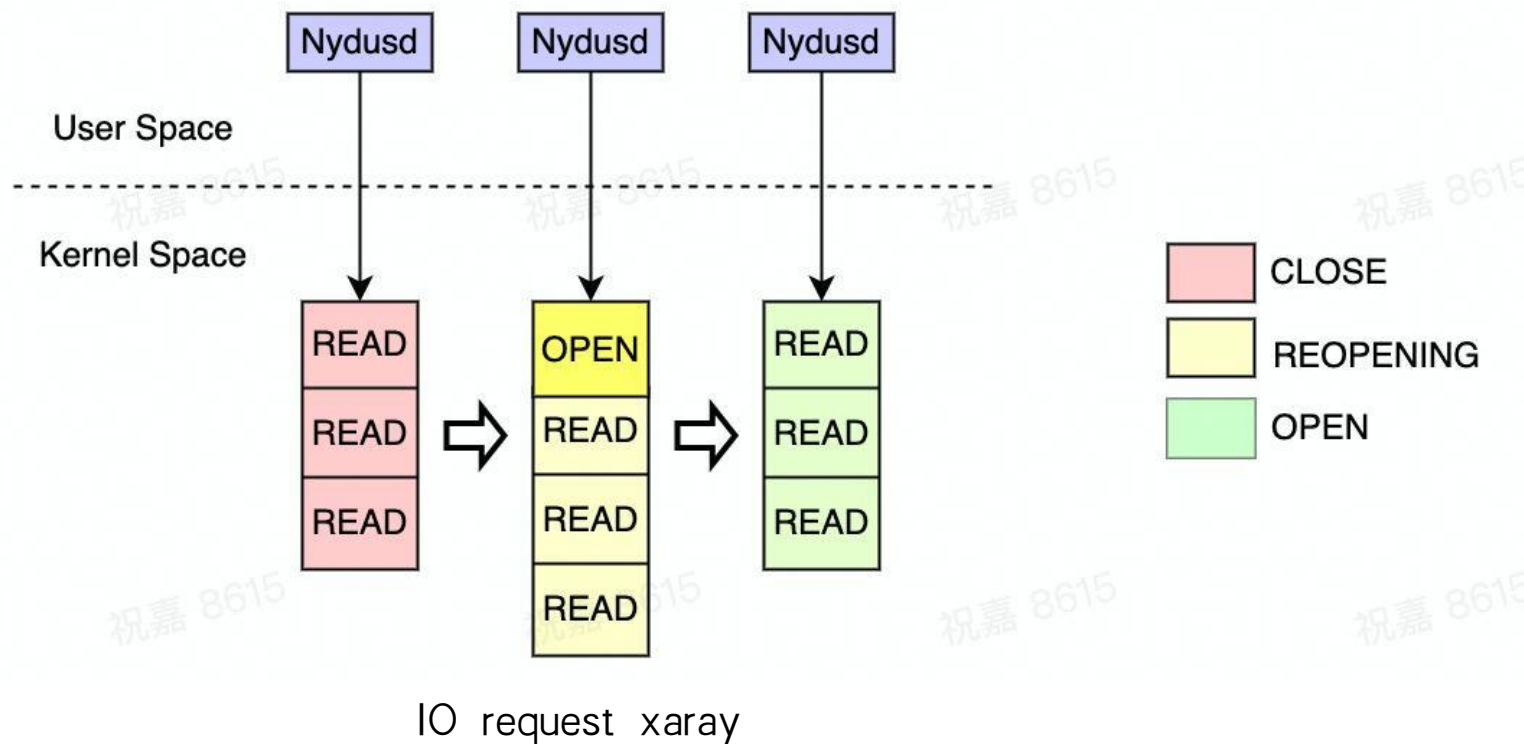
- /dev/cachefiles fd: 用户态使用UDS保存和恢复



## 稳定性提升：User Daemon Failover

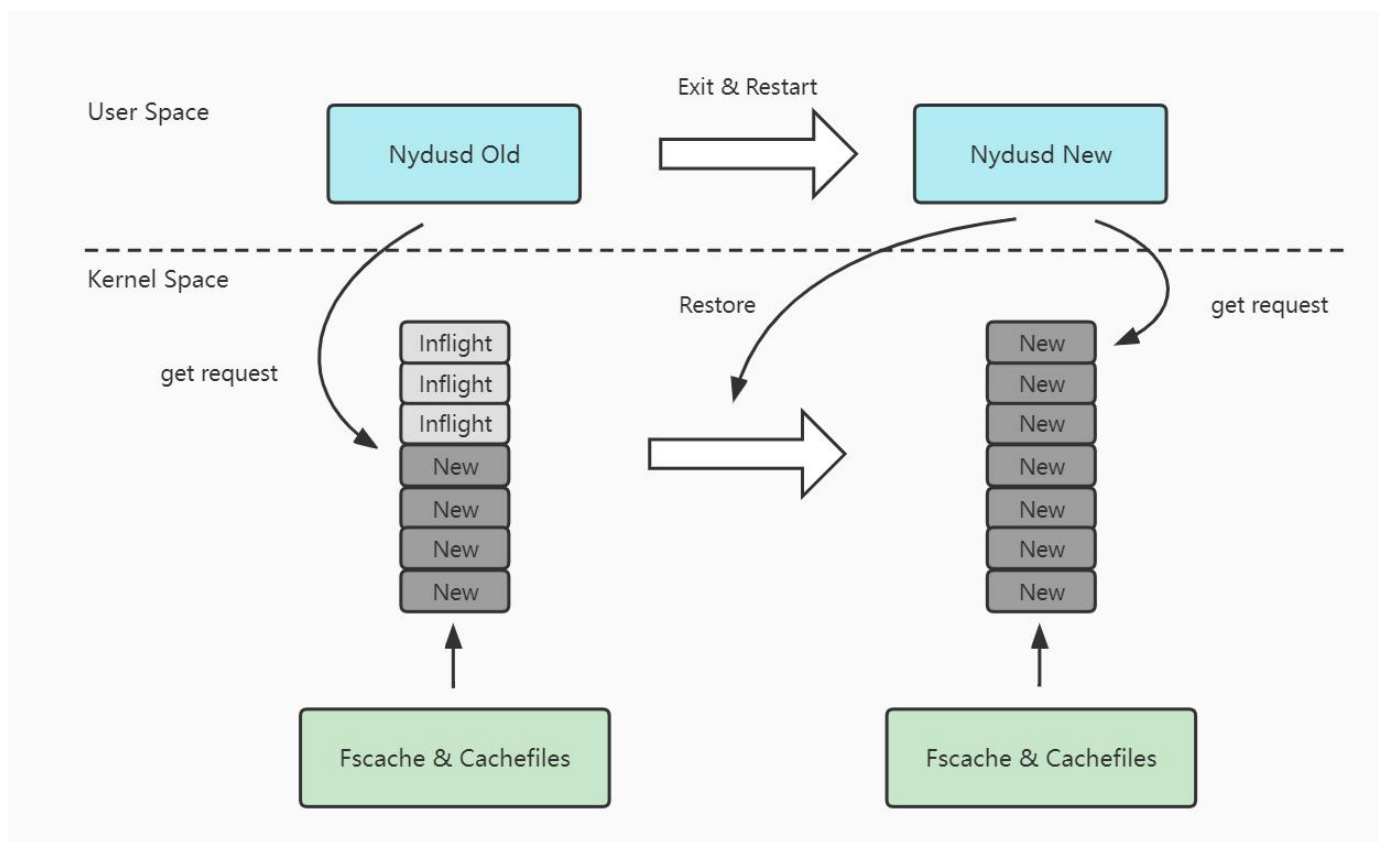
fd资源保持与恢复

- 缓存文件匿名fd：由按需IO触发，重新生成OPEN请求打开文件



# 稳定性提升：User Daemon Failover

- 恢复IO请求



# 稳定性提升：User Daemon Failover

- 进展：正在推动进入upstream

- <https://lore.kernel.org/all/20221014080559.42108-1-zhujia.zj@bytedance.com/>



## 实践与探索——New Features

01

Failover

02

**EROFS Shared  
Domain**

03

Daemonless

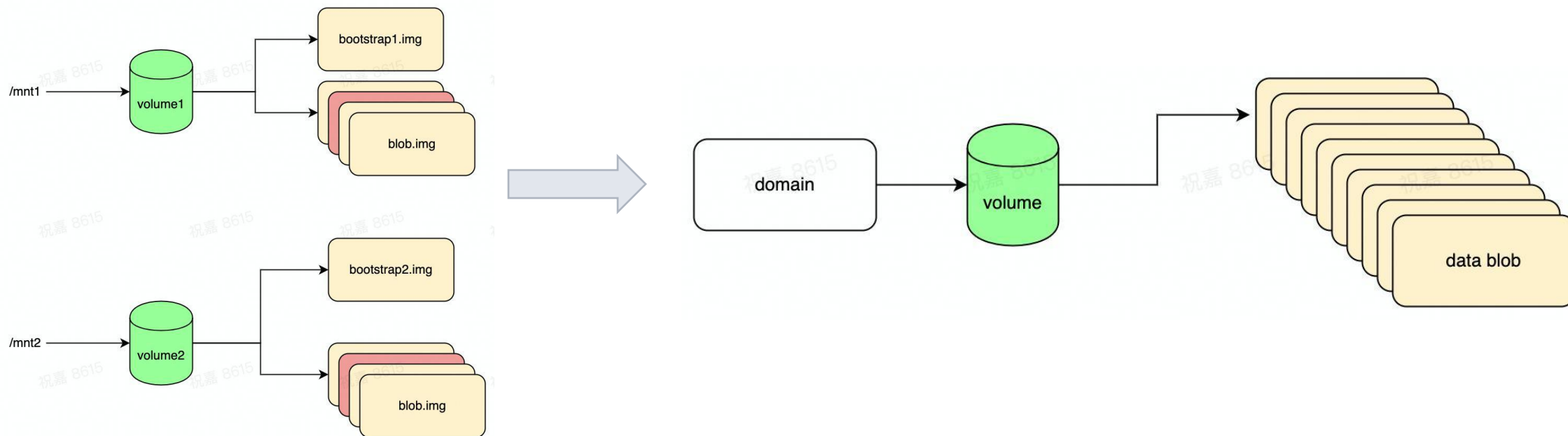
# 数据去重：EROFS Shared Domain

背景：

- 不同EROFS挂载点间，相同的镜像无法共享

预期目标：

- 实现镜像layer粒度的共享



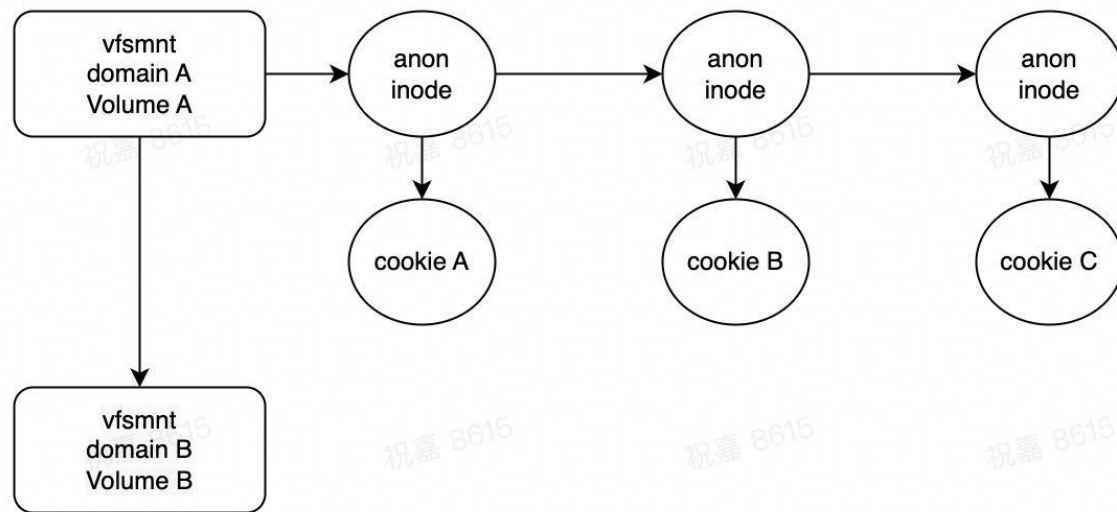
# 数据去重：EROFS Shared Domain

## 设计

- EROFS侧抽象Shared Domain模型，对应fscache volume
- pseudo挂载点
  - 管理domain生命周期
  - anonymous inodes 管理镜像文件生命周期，对应fscache cookie (blob)

## 使用

- EROFS挂载选项'domain\_id'
- 创建/加共享域：同一共享域内所有挂载点共享镜像
- 不使用domain\_id：不与任何挂载点共享镜像数据



# 数据去重：EROFS Shared Domain

CACHE	REF	VOLS	OBJS	ACCES	S	NAME
=====	=====	=====	=====	=====	=====	=====
00000009	3	2	11	1	A	CacheFiles
VOLUME	REF	nCOOK	ACC	FL	CACHE	KEY
=====	=====	=====	=====	=====	=====	=====
00000009	7	6	1	00	CacheFiles	erofs,7ea193b433cad62d1b909ecd6a948201fb17ff5d2aca3284154bad090109b186
0000000a	6	5	1	00	CacheFiles	erofs,32e72ef854d48d4850f6d677b5a085dc25cd69250129636ae0cd8322f7722201
COOKIE	VOLUME	REF	ACT	ACC	S	FL DEF
=====	=====	=====	=====	=====	=====	=====
00000022	00000009	2	1	0	A	6134 37656131393362343333636164363264316239303965636436613934383230316662313766663564326163613332383431353462616430393031303962313836
00000023	00000009	2	1	0	A	6134 35313265306664646564613130653366333763613563373031373261653931326134616664323133366337393238373433653861346636303566373138656264
00000024	00000009	2	1	0	A	6134 63373333386136343862316332636362353364613737356166643462323463346663666133353839353735383462303566326339343833363463306361356662
00000025	00000009	2	1	0	A	6134 63623834313834613633643962336234333936633961626266656131323066643661356537353132663933626431326339393464663761343063303937623535
00000026	00000009	2	1	0	A	6134 35393733316532336265333761616335336336646564313839343464663434313536616131343131326234343639643632343231666164646336393937343532
00000027	00000009	2	1	0	A	6134 32363432346138323633356431353761623332376661356236373764633131646536353563666132666338303839656630323735643333346436643430623237
00000028	0000000a	2	1	0	A	6134 33326537326566383534643438643438353066366436373762356130383564633235636436393235303132393633366165306364383332326637373232323031
00000029	0000000a	2	1	0	A	6134 35313265306664646564613130653366333763613563373031373261653931326134616664323133366337393238373433653861346636303566373138656264
0000002a	0000000a	2	1	0	A	6134 63373333386136343862316332636362353364613737356166643462323463346663666133353839353735383462303566326339343833363463306361356662
0000002b	0000000a	2	1	0	A	6134 63623834313834613633643962336234333936633961626266656131323066643661356537353132663933626431326339393464663761343063303937623535
0000002c	0000000a	2	1	0	A	6134 35393733316532336265333761616335336336646564313839343464663434313536616131343131326234343639643632343231666164646336393937343532

CACHE	REF	VOLS	OBJS	ACCES	S	NAME
=====	=====	=====	=====	=====	=====	=====
00000007	2	1	7	1	A	CacheFiles
VOLUME	REF	nCOOK	ACC	FL	CACHE	KEY
=====	=====	=====	=====	=====	=====	=====
00000007	8	7	1	00	CacheFiles	erofs,ABigDomain
COOKIE	VOLUME	REF	ACT	ACC	S	FL DEF
=====	=====	=====	=====	=====	=====	=====
0000001a	00000007	2	1	0	A	6134 37656131393362343333636164363264316239303965636436613934383230316662313766663564326163613332383431353462616430393031303962313836
0000001b	00000007	2	1	0	A	6134 35313265306664646564613130653366333763613563373031373261653931326134616664323133366337393238373433653861346636303566373138656264
0000001c	00000007	2	1	0	A	6134 63373333386136343862316332636362353364613737356166643462323463346663666133353839353735383462303566326339343833363463306361356662
0000001d	00000007	2	1	0	A	6134 63623834313834613633643962336234333936633961626266656131323066643661356537353132663933626431326339393464663761343063303937623535
0000001e	00000007	2	1	0	A	6134 35393733316532336265333761616335336336646564313839343464663434313536616131343131326234343639643632343231666164646336393937343532
0000001f	00000007	2	1	0	A	6134 32363432346138323633356431353761623332376661356236373764633131646536353563666132666338303839656630323735643333346436643430623237
00000020	00000007	2	1	0	A	6134 33326537326566383534643438643438353066366436373762356130383564633235636436393235303132393633366165306364383332326637373232323031

Shared Domain



## 数据去重：EROFS Shared Domain

### • 进展

- 已进入upstream v6.1-rc1
- <https://lore.kernel.org/all/20220916085940.89392-1-zhujia.zj@bytedance.com/>

### • 未来工作

- 发挥内核文件系统优势，探索page cache sharing

## 实践与探索——New Features

01

Failover

02

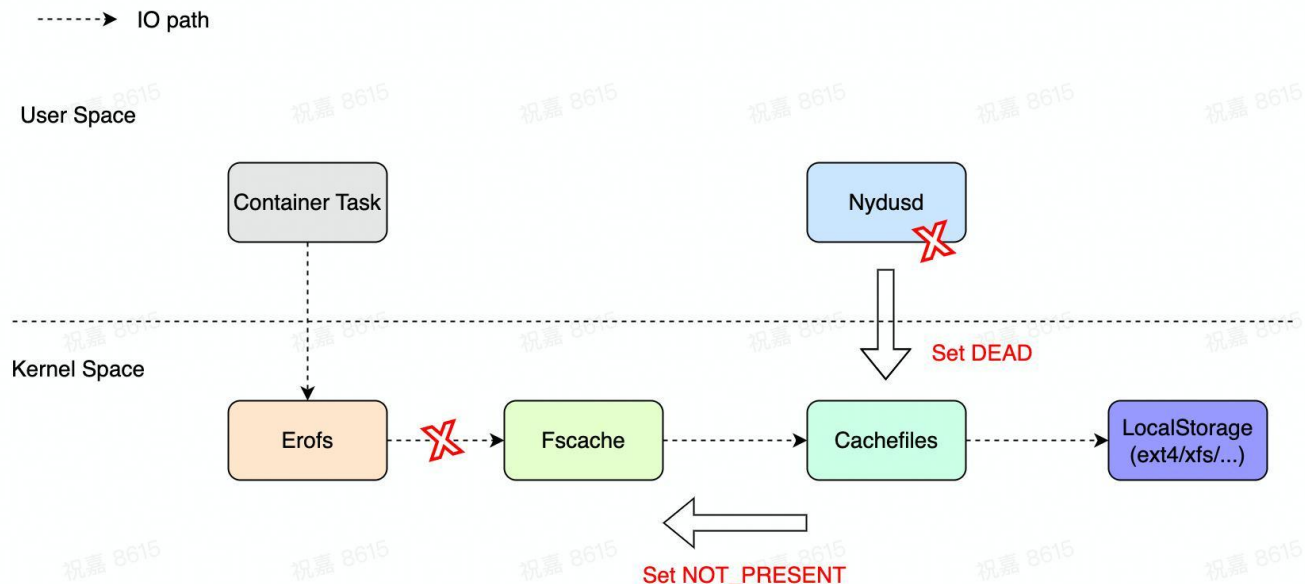
EROFS Shared  
Domain

03

**Daemonless**

## 0开销：Daemonless

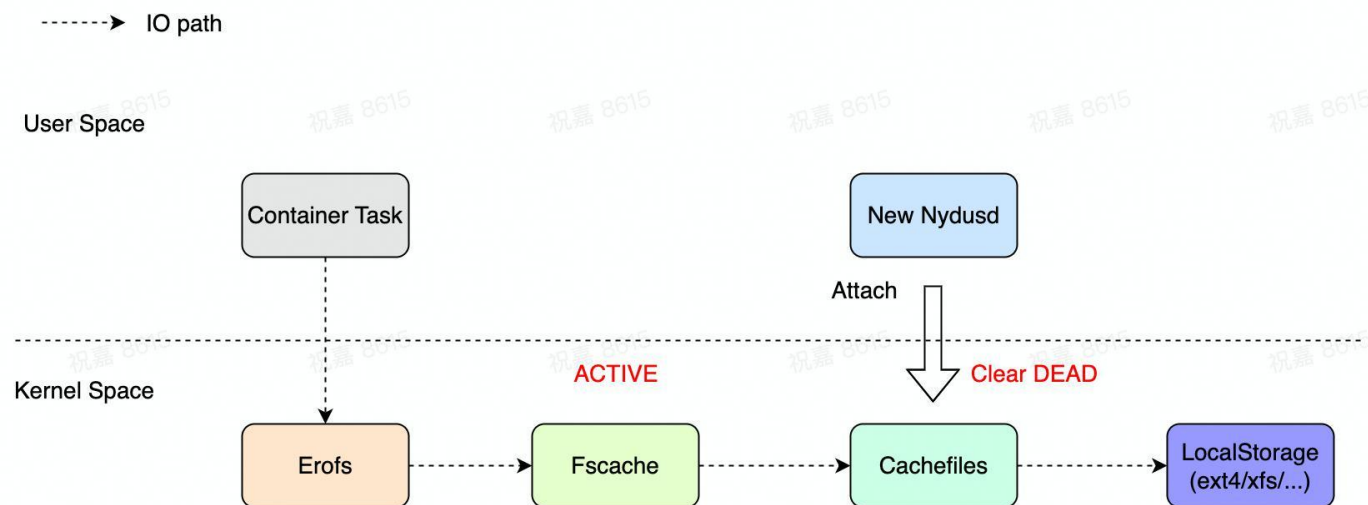
- 背景**：按需加载完成后，IO转为纯内核路径，不再需要用户进程参与，可使其退出，节省资源
- 面临的挑战**：用户进程退出后，fscache相关内核结构被置为DEAD状态，EROFS无法从fscache中读取数据
- 设计目标**：cachefiles实现新运行模式，解除本地数据访问对用户daemon的依赖



# 0开销：Daemonless

## 设计目标

- daemon 退出后，内核不释放相关结构，EROFS依然可以访问本地缓存文件
- daemon启动，可重新接管上述本地缓存文件，继续进行cache维护管理。





# 0开销：Daemonless

进展

代码已验证完成，即将提交内核社区

# Thanks