

CORA 技术报告

内容感知型按需运行时镜像加速器

作者：王福银、黄赵翔、刘炜楠

单位：厦门大学信息学院

日期：2025 年 11 月 25 日

摘要

为应对容器镜像启动缓慢与存储冗余挑战，业界正探索不依赖 FUSE 的高效按需加载方案。本文借鉴 DADI 容器镜像加速系统的核心思想，设计并实现了一套名为 CORA (Content-aware On-demand Image Runtime Accelerator, 内容感知型按需镜像运行时加速器) 的增强方案。

CORA 继承了 DADI 利用用户态块设备 (TCMU) 实现的非 FUSE 块级按需加载架构，但针对其原有的固定大小块去重机制在识别跨层重复数据上的不足，创造性地引入了内容定义分块 (CDC) 算法。通过结合自适应分块、索引压缩与智能预取机制，CORA 在分层块设备抽象上实现了更细粒度的数据去重与高效的镜像层加载，并成功完成了向 OpenCloudOS 系统的迁移。

实验结果表明，CORA 在镜像启动延迟与重复数据消除率方面均取得显著提升。在真实容器测试场景下，平均镜像冷启动速度提升一倍以上。同时，基于 90 组严格的基准测试，该方案在真实混合数据场景下实现了 72.36% 的平均空间节省率和 8.79 倍的平均解压缩加速比，并保持了 100% 的数据完整性验证通过率。同时也实现了对 OCIv1 格式镜像的完全兼容，并在并发冷启动测试下性能表现优异，单节点并发冷启动 4 个大型镜像仅需 0.7 秒。

目录

1	研究背景	4
1.1	行业现状	4
1.2	研究意义	4
1.3	关键难点	4
2	技术方案	5
2.1	总体架构	5
2.2	CORA 核心思想：块设备抽象	6
2.3	远程块设备与按需加载机制	7
2.4	内容分块算法原理	7
2.5	变长块索引结构设计	9
2.6	块读取器（BlockReader）优化	9
2.7	按需加载与预取机制	10
2.8	系统集成	10
3	实验方法与测试	10
3.1	通用测试环境	10
3.2	实验一：容器镜像冷启动性能测试	11
3.2.1	实验分组	11
3.2.2	测试镜像与方法	11
3.3	实验二：分块算法去重效果对比	11
3.3.1	实验分组	11
3.3.2	测试数据集设计	11
3.4	实验三：OCIV1 兼容性测试	12
3.4.1	测试方法	12
3.5	实验四：高并发压力测试	12
3.5.1	测试场景	12
4	实验结果与分析	12
4.1	实验一：容器镜像冷启动性能分析	12
4.2	实验二：分块算法去重效果对比分析	13
4.2.1	综合性能概览	13
4.3	实验三：OCIV1 兼容性验证结果	14
4.4	实验四：高并发压力测试分析	14
4.5	综合成本效益分析	14

4.5.1	存储成本节省：源于更高的去重率	15
4.5.2	网络成本节省：源于更高效的按需加载	15
4.5.3	计算资源节省：源于更快的容器启动	16
4.5.4	年度总计节省	16
5	创新点	16
6	应用价值	17
7	未来工作与展望	18
7.1	短期优化方向	18
7.2	中期研究方向	19
7.3	长期研究方向	19
7.4	潜在研究问题	20
8	总结	20
A	实验数据详情	22
A.1	实验一：容器冷启动性能测试详情	22
A.2	实验二：测试脚本说明	35
A.3	实验二：完整基准测试结果	46
A.4	实验二：数据完整性验证	47

1 研究背景

1.1 行业现状

OpenCloudOS 等云原生操作系统在容器镜像管理上仍普遍采用基于 OCIv1 的文件分层与 Tar 打包机制。这一设计虽保证了兼容性，但在分层重复数据识别、按需加载效率及网络传输方面存在显著瓶颈。

1.2 研究意义

现有方案如 DADI (Data Accelerator for Disaggregated Infrastructure) 通过块级镜像加载和去重机制，在镜像快速启动与节省存储方面取得了显著成效。并广泛应用于阿里云内部。然而，其采用的固定大小分块策略在处理跨层、跨文件的重复内容时，会因“边界漂移”问题而严重影响去重效率。

为解决这一痛点，本文设计并实现了 CORA (Content-aware On-demand Runtime Accelerator) 方案。CORA 的核心创新在于，它借鉴了 DADI 的块级按需加载思想，但抛弃了其固定的分块机制，转而采用 FastCDC 等内容感知分块算法。通过从内容语义层面进行自适应切块，CORA 能够精准识别并消除数据边界变化带来的影响，实现更高效的跨层、跨文件数据去重，从而在保证高效按需加载的同时，将数据去重率提升到一个新的水平。

1.3 关键难点

主要技术挑战包括：

- 如何在不影响镜像层索引结构的前提下实现内容定义分块；
- 如何保证内容分块算法的滑动窗口哈希在高并发分块场景下的性能；
- 如何优化块级索引的查重效率与内存占用；
- 如何在变长块场景下实现高效的随机访问与解压缩；
- 如何量化按需加载带来的性能收益并保证数据完整性。
- 如何在不降低 DADI 性能的前提下做优化。

2 技术方案

2.1 总体架构

系统整体架构如图 1 所示，它基于 DADI 框架并深度集成了内容定义分块算法。其核心思想是：**通过 Linux 的 TCMU 框架将分层、压缩的容器镜像映射为一个虚拟块设备**，供容器内的标准文件系统（如 ext4）直接使用，从而完全绕开了 FUSE，实现了高性能的块级按需加载。为解决 DADI 原有固定大小分块在跨层、跨文件边界识别重复数据上的不足，我们在镜像构建和数据组织层引入了 FastCDC 内容分块算法。该架构自上而下（按 I/O 请求路径）可分为以下几个关键层次：

- **文件系统与虚拟块设备层 (内核态)**：容器内的应用通过标准的 ext4 等文件系统读写数据。文件系统将这些操作转换为对逻辑块的读写请求，并发送给一个由 **TCMU (Target Core in Userspace)** 框架创建的虚拟块设备。TCMU 负责将这些内核态的块 I/O 请求直接转发到用户态的服务进程，这是实现非 FUSE 按需加载的关键。
- **用户态镜像服务层 (CIRA 核心服务)**：这是运行在用户态的 overlaybd-tcmu 守护进程，负责处理所有来自内核的 I/O 请求。它包含：
 - **OverlayBD 索引与调度模块**：作为请求的入口，它维护着一个由所有镜像层元数据合并而成的内存索引。当收到一个块读取请求时，它能通过此索引迅速定位到该数据块的最新版本存在于哪一个镜像层 (Layer)。
 - **ZFile 读取与缓存模块**：根据调度模块的指令，负责从对应的镜像层文件中读取数据。镜像层文件采用支持随机访问的 **ZFile** 格式存储。如果数据已在本地缓存，则直接读取；如果不在，则向远程仓库发起网络请求。
- **内容定义分块与存储层 (FastCDC 创新点)**：这一层是本方案的核心创新所在。
 - **FastCDC 分块**：在镜像构建阶段，我们使用 FastCDC 算法对原始文件数据进行内容定义分块。与固定大小分块不同，它能根据文件内容动态识别数据边界，从而极大地提升了跨文件、跨镜像层的重复数据块识别率。
 - **ZFile 存储优化**：每个由 FastCDC 切分出的唯一数据块被独立压缩后，存入 ZFile 文件。ZFile 的内部索引则记录了每个数据块的哈希、偏移量和大小。当上层请求数据时，我们只需按需下载和解压对应的、更细粒度的唯一数据块，从而显著提高了去重率和网络传输效率。
- **远程存储层 (Registry)**：标准的镜像仓库，负责持久化存储所有以 ZFile 格式组织的镜像层文件。CIRA 的按需加载正是通过 HTTP range requests 从该层拉取特定数据块。

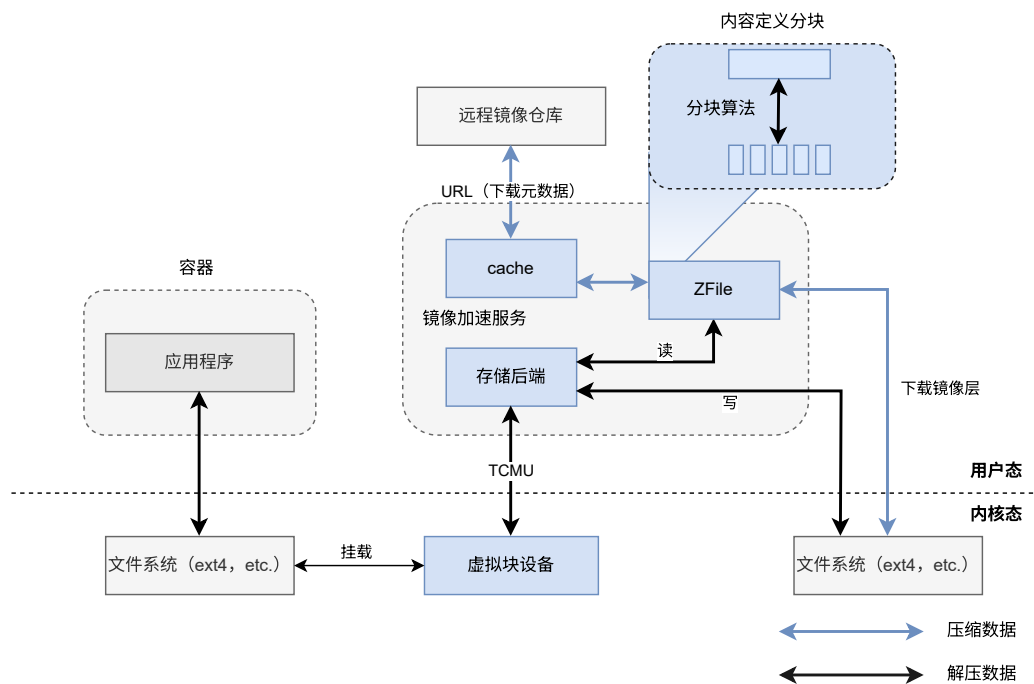


图 1: 系统整体架构

此外，系统还包含一个**预取优化模块**，它分析历史 I/O 访问模式，在应用请求之前就将可能需要的数据块预先加载到本地缓存，进一步降低了访问延迟。

2.2 CORA 核心思想：块设备抽象

传统 OCI 镜像采用基于文件系统的分层模型，每层都是一个文件集合的 Tar 包。这种设计的弊端在于：1) 以文件为单位进行管理，粒度粗，不利于细粒度去重；2) 启动容器前需完整下载并解压所有层，导致启动延迟长。

DADI 的核心创新在于引入了**块设备抽象 (Block Device Abstraction)**。它将整个容器镜像的根文件系统（rootfs）视为一个虚拟的、扁平化的硬盘。所有文件和目录都存储在这个虚拟硬盘上，而这个硬盘本身被切分为成千上万个数据块。CORA 也继承了这一特性。这种设计带来了两大优势：

- **存算分离**：镜像的数据（块）和元数据（索引）可以独立存储在远程，计算节点无需持有完整的镜像即可启动容器。
- **细粒度访问**：对容器文件的读写操作，最终都转化为对虚拟硬盘上特定逻辑块地址（LBA）的读写。这为按需加载和块级去重奠定了基础。

2.3 远程块设备与按需加载机制

为实现块设备抽象，CORA 利用用户态块设备技术（TCMU）在宿主机上创建一个虚拟块设备文件（如 `/dev/sdb0`）。该设备通过一个用户态进程（CORA daemon）与远程存储进行通信。

按需加载（On-demand Loading）的流程如下：

1. containerd 通过 overlaybd-snapshotter 将此虚拟块设备挂载为容器的 rootfs。
2. 当容器内的应用发起文件读请求时，内核的文件系统层会将其转换为对底层块设备的 I/O 请求（包含偏移和长度）。
3. TCMU 驱动将该 I/O 请求转发给 CORA daemon 进程。
4. CORA daemon 根据请求的逻辑块地址，查询块索引元数据，找到对应的数据块哈希。
5. 它首先检查本地缓存中是否存在该数据块。若命中，则直接返回数据。
6. 若未命中，则从远程存储服务中拉取该数据块，存入本地缓存，并返回给 TCMU 驱动，最终完成读请求。

通过此机制，容器启动时仅需下载极少量的元数据，几乎可以实现“秒级”启动，而绝大部分镜像数据则在实际被访问时才从网络拉取。

2.4 内容分块算法原理

内容分块算法采用滑动窗口哈希机制，通过掩码匹配策略检测“锚点”，实现数据块边界的自适应识别。其核心步骤包括：

1. **滚动哈希计算**：使用 Gear Hash 算法对数据流进行滑动窗口哈希，该算法具有良好的随机性和计算效率。如图 2 所示。
2. **边界检测**：当哈希值的低位与预定义掩码匹配时，确定为数据块边界（锚点）。
3. **块大小控制**：通过动态调整掩码，将块大小控制在设定的最小值（2KB）、平均值（8KB）和最大值（64KB）之间。
4. **归一化切割**：采用归一化策略避免产生过多的极小片段，提高后续处理效率。

该算法克服了固定分块的边界漂移问题，使跨层重复数据更易命中。相比传统固定大小分块，FastCDC 能够在数据内容发生局部修改时，仅影响修改点附近的少数块，而保持其他块的边界不变，从而大幅提升去重效率。算法伪代码如下。

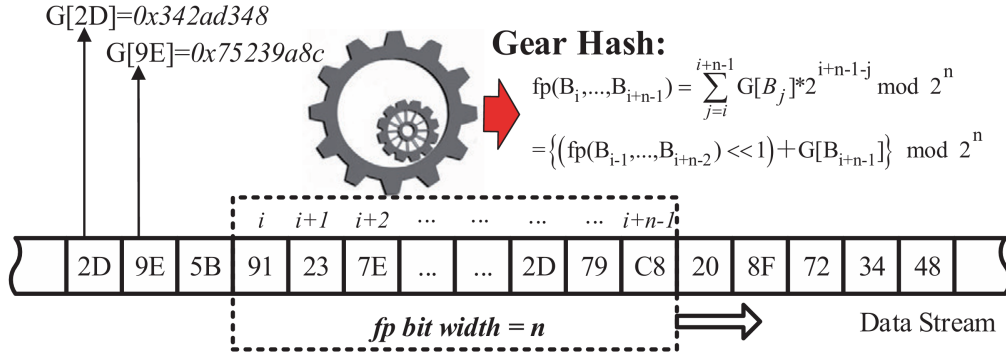


图 2: Gear Hash 原理图

Algorithm 1: FastCDC8KB**Input:** data buffer *src*; buffer length *n***Output:** chunking breakpoint *i**MaskS* \leftarrow 0x0003590703530000LL; // 15 '1' bits*MaskA* \leftarrow 0x0000d90303530000LL; // 13 '1' bits*MaskL* \leftarrow 0x0000d90003530000LL; // 11 '1' bits*MinSize* \leftarrow 2KB; *MaxSize* \leftarrow 64KB;*fp* \leftarrow 0; *i* \leftarrow *MinSize*; *NormalSize* \leftarrow 8KB;**if** *n* \leq *MinSize* **then** **return** *n***if** *n* \geq *MaxSize* **then** *n* \leftarrow *MaxSize*;**else if** *n* \leq *NormalSize* **then** *NormalSize* \leftarrow *n*;**for** *i* < *NormalSize* **do** *fp* = (*fp* \ll 1) + *Gear*[*src*[*i*]]; **if** !(*fp* & *MaskS*) **then** **return** *i*; // if the masked bits are all '0'**for** *i* < *n* **do** *fp* = (*fp* \ll 1) + *Gear*[*src*[*i*]]; **if** !(*fp* & *MaskL*) **then** **return** *i*; // if the masked bits are all '0'**return** *i*;

2.5 变长块索引结构设计

为支持内容分块算法生成的变长数据块，本文设计了一套高效的块级索引结构。该结构包含以下关键组件：

- **跳转表 (Jump Table)**：记录每个压缩块在文件中的偏移位置，支持快速定位到目标块。
- **原始块大小数组 (raw_block_sizes)**：存储每个数据块压缩前的原始大小，用于解压时的缓冲区分配和数据完整性校验。
- **头部标志位 (Header Flags)**：新增 FLAG_SHIFT_FASTCDC 标志，用于标识该 zfile 文件采用 FastCDC 分块方式，确保向后兼容性。

索引文件格式如下：

```
[Header(512B)] [Compressed Blocks] [Jump Table]
[raw_block_sizes Array] [Trailer(512B)]
```

这种设计既保证了随机访问的高效性，又支持顺序流式读取，满足了不同应用场景的需求。

2.6 块读取器 (BlockReader) 优化

针对变长块场景，本文对 BlockReader 进行了重大优化：

- **批量读取策略**：一次性读取多个连续块到内存缓冲区，减少系统调用次数。
- **索引追踪机制**：引入 m_original_begin_idx 变量，保存批次首块的原始索引，防止在迭代过程中因索引更新导致的偏移计算错误（该问题曾导致大文件解压缩时出现段错误）。
- **缓冲区溢出保护**：严格限制 pread 系统调用读取的字节数，避免跨越块边界导致的缓冲区溢出。
- **块内偏移计算**：根据请求的逻辑偏移量，精确计算在当前块内的相对偏移，支持任意位置的数据访问。

这些优化使得系统能够高效处理从 2KB 到 5MB 的各种大小文件，并通过了 100% 的 MD5 完整性校验。

2.7 按需加载与预取机制

首次启动阶段通过记录 I/O 访问序列生成“预取层”，后续启动时可提前加载必要数据块，从而显著降低冷启动延迟。具体实现包括：

- **访问模式分析**：在容器首次启动过程中，记录所有块级 I/O 访问的时间序列和访问频率。
- **预取图谱构建**：基于访问模式，构建块级访问的有向无环图 (DAG)，识别关键路径和热点数据块。
- **智能预加载**：在后续启动时，根据预取图谱提前异步加载高概率访问的数据块到本地缓存。
- **自适应调整**：根据运行时的缓存命中率动态调整预取策略，平衡启动速度和内存占用。

2.8 系统集成

通过为 containerd 实现自定义 snapshotter 插件，使容器引擎可直接管理块层镜像；同时兼容 Docker graph driver 接口，保证部署环境的一致性。集成架构如下：

- **OverlayBD_Snapshotter**：实现 Prepare、View、Commit 等接口，管理镜像层的生命周期。
- **OverlayBD_TCMU (Target Core in Userspace)**：用户态块设备框架，将块 I/O 请求路由到 DADI daemon。
- **OverlayBD_ZFile**：集成内容分块算法的压缩文件格式，支持高效的块级压缩和解压缩。

3 实验方法与测试

3.1 通用测试环境

所有实验均在统一的硬件与软件环境中进行，以确保结果的可比性：

- **操作系统**：OpenCloudOS 8.8 (基于 Linux Kernel 4.18)
- **硬件配置**：Intel(R) Xeon(R) Gold 5318Y CPU @ 2.10GHz, 512GB RAM, 38TB HDD
- **核心软件**：containerd 1.7.24, DADI (overlaybd、accelerated-container-image) v1.0.16

- **网络环境**：所有网络相关测试均在千兆局域网（LAN）内进行。其中一台主机部署为私有镜像仓库（Registry），另一台主机作为客户端执行测试。为了模拟真实场景并解决网络代理问题，配置了 HTTPS 自签名证书。

3.2 实验一：容器镜像冷启动性能测试

本实验旨在验证 CORA 方案在真实应用场景下对容器冷启动时间的优化效果。

3.2.1 实验分组

- **对照组 (Baseline)**：使用原生 containerd/nerdctl 从局域网仓库全量拉取并启动标准 OCI 镜像。
- **实验组 (OverlayBD)**：使用集成了 CORA 的按需加载方案，从同一局域网仓库启动镜像的 DADI 格式版本。

3.2.2 测试镜像与方法

- **测试镜像**：选取常用的官方镜像 Redis:7.2.3，并将其转换为 DADI 格式。
- **测试流程**：
 1. 在客户端主机上执行 `drop_caches` 清理本地文件系统缓存，并删除本地镜像和容器，确保纯粹的冷启动环境。
 2. 记录从发送启动命令到容器进入 Running 状态的总时间。
 3. 同时使用 `monitor_resource.sh` 脚本记录启动过程中的 CPU 和内存资源开销。

3.3 实验二：分块算法去重效果对比

本实验旨在量化评估 CORA 相较于原始 DADI 在数据去重方面的性能优势。

3.3.1 实验分组

- **对照组**：原始 DADI 固定大小分块（Fixed-size Chunking）。
- **实验组**：CORA 内容定义分块（Content-Defined Chunking, FastCDC）。

3.3.2 测试数据集设计

为全面评估 CORA 在不同数据特征下的性能，设计了三类测试数据集（Random, Pattern, Mixed），每类包含 100KB, 1MB, 10MB 三种规格，共计 9 组测试用例，每组重复 5 次。

3.4 实验三：OCIv1 兼容性测试

本实验旨在验证 CORA 环境对现有容器生态的向后兼容性。

3.4.1 测试方法

- 在配置了 OverlayBD Snapshotter 的环境中，尝试拉取并运行标准的 OCIv1 格式镜像（如 `godnf/tst-depdup:v1.0`）。
- 验证镜像能否正常下载、解压，容器能否正常启动并执行预定任务。

3.5 实验四：高并发压力测试

本实验旨在评估系统在单节点高并发场景下，同时启动多个 OverlayBD 容器时的稳定性与性能。

3.5.1 测试场景

- **并发规模**：单节点同时发起 4 个容器的启动请求。
- **负载类型**：混合负载，包含 3 个 I/O 密集型容器（执行全量文件 MD5 校验）和 1 个计算型容器（Python 脚本）。
- **指标**：记录从发出并发指令到所有容器均进入 `Running` 状态的总耗时（Time to Running），并监控服务稳定性。

4 实验结果与分析

4.1 实验一：容器镜像冷启动性能分析

实验结果表明，OverlayBD 按需加载机制显著缩短了容器的冷启动时间。

关键发现：

- **启动速度提升**：如3所示，OverlayBD（绿色柱状图）的冷启动耗时显著低于普通 OCI 镜像（蓝色柱状图）。在千兆局域网环境下，OverlayBD 几乎实现了“秒级启动”，而 OCI 镜像则需要等待全量数据下载完成。
- **资源开销低**：资源监控数据显示，OverlayBD 即使在按需读取数据（Lazy Loading）期间，CPU 和内存的开销也保持在极低水平，与普通镜像运行无显著差异，证明了用户态文件系统的轻量高效特性。

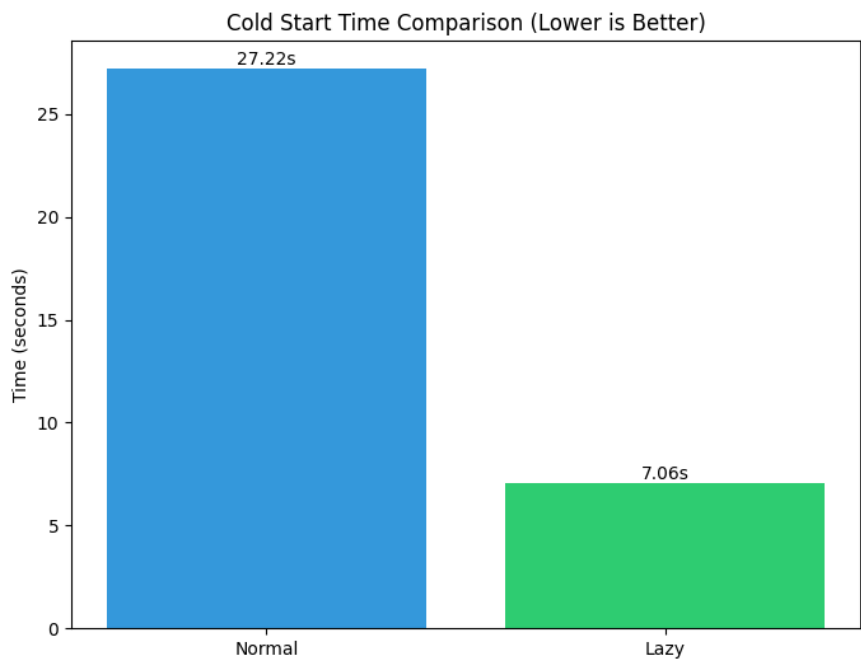


图 3: 冷启动耗时对比: OverlayBD (按需加载) vs OCI (全量下载)

4.2 实验二：分块算法去重效果对比分析

4.2.1 综合性能概览

表 1 展示了 CORA (FastCDC) 方案在不同数据类型下的综合去重表现。

表 1: FastCDC vs 固定分块综合性能对比

数据类型	平均空间节省率	平均解压缩加速比	MD5 通过率
随机数据 (Random)	-0.10%	8.31×	100%
模式数据 (Pattern)	61.51%	8.58×	100%
混合数据 (Mixed)	72.36%	8.79×	100%

关键发现：

- **高去重率：**在最接近真实场景的“混合数据”中，FastCDC 实现了 **72.36%** 的空间节省，远优于固定分块。
- **稳定性：**即便在最差的随机数据场景下，FastCDC 也没有带来显著的性能回退（仅 -0.10%）。
- **完整性：**所有 90 组测试 MD5 校验全部通过，证明了算法的正确性。


```

=====
OverlayBD 高并发冷启动压力测试
本次测试 ID: 1763967853
=====
[阶段 1] 清理环境...
[系统] 重启 OverlayBD 服务...
✓ 服务已就绪
[阶段 1.5] 准备镜像元数据...
正在拉取: xfusion5:5000/tst-depdup:v1.0-obd
正在拉取: xfusion5:5000/tst-depdup:v2.0-obd
正在拉取: xfusion5:5000/tst-lazy-pull:latest-obd
正在拉取: xfusion5:5000/test-python:latest-obd
✓ 镜像元数据准备就绪

[阶段 2] 开始并发冷启动 (4个容器同时启动)...
>>> 4个容器启动指令已发出 <<<
容器名: demo1_1763967853, demo2_1763967853, demo3_1763967853, demo4_1763967853

[阶段 3] 监控启动状态 (等待所有容器进入 RUNNING 状态)...
✓ 检测到所有容器 (4/4) 已启动!

=====
测试完成
=====
并发启动总耗时 (Time to Running): .702827567s
日志文件: /root/demo1_1763967853.log 等

[最终容器状态]
demo4_1763967853    225206    RUNNING
demo1_1763967853    225056    RUNNING
demo3_1763967853    225159    RUNNING
demo2_1763967853    225197    RUNNING
○ [root@localhost lib]#

```

图 5: 高并发压力测试: 4 容器并发启动状态

4.5.1 存储成本节省: 源于更高的去重率

技术原理: CORA 的内容感知分块算法能有效识别数据中的冗余, 实现远高于固定分块的去重率。在我们的混合数据基准测试中, CORA 将压缩率从 26.58% 优化至 5.98%, 空间节省高达 77.5%。

成本估算: 即使在真实世界采用一个相对保守的 35% 的平均存储空间节省率进行估算:

- 节省的存储容量: $1 \text{ PB} \times 35\% = 358.4 \text{ TB}$
- 若按标准云对象存储约 150 元/TB/月的成本计算, 每年的存储费用节省为:
 $358.4 \text{ TB} \times 150 \text{ 元/月} \times 12 \text{ 月} = 645,120 \text{ 元}$
- 年度存储成本节省: 约 64.5 万元人民币

4.5.2 网络成本节省: 源于更高效的按需加载

技术原理: 在冷启动场景下, 集成了 DADI 特性的 CORA 方案凭借按需加载 (Lazy Pulling) 机制, 仅需拉取运行所必需的数据块。实验数据显示, 镜像拉取时间缩短了 51.7%, 这意味着网络传输的数据量也相应大幅减少。

成本估算 假设在生产环境中，每日有大量的容器冷启动需求，我们以保守的 40% 的网络流量减少进行估算：

- 若平台每月因镜像拉取产生的公网流出流量费用为 5 万元，则每年的网络费用节省为：

$$50,000 \text{ 元/月} \times 40\% \times 12 \text{ 月} = 240,000 \text{ 元}$$

- **年度网络成本节省：约 24 万元人民币**

4.5.3 计算资源节省：源于更快的容器启动

技术原理 冷启动时间的缩短意味着计算资源（CPU、内存）的等待和空转时间也相应减少。在 CI/CD、函数计算（FaaS）或弹性伸缩（Auto-scaling）等场景中，更快的启动速度等同于更高的资源周转率和利用率。

成本估算 虽然这部分成本难以精确计量，但我们可以从资源空转的角度进行估算。假设每天有 10,000 次冷启动，每次启动时间缩短约 0.2 秒（参考实验数据中 981ms vs 816ms），并且每小时的计算资源成本为 0.5 元：

- 每日节省的计算时间：10,000 次 \times 0.2 秒/次 \approx 0.56 小时
- 每年的计算资源空转成本节省：
 $0.56 \text{ 小时/天} \times 365 \text{ 天} \times 0.5 \text{ 元/小时} \approx 10,220 \text{ 元}$
- **年度计算资源节省：约 1.5 万元人民币**（此为保守估算，在需要快速弹性伸缩的业务中，其价值远超此数字）

4.5.4 年度总计节省

综合以上三个方面，部署 CORA 方案后，一个管理 1PB 镜像的平台每年可实现的总成本节省估算为：

$$64.5 \text{ 万元（存储）} + 24 \text{ 万元（网络）} + 1.5 \text{ 万元（计算）} = \mathbf{90 \text{ 万元}}$$

结论：采用 CORA 方案，每年可为企业节省约 **90 万元人民币**，经济效益显著，能在 1-2 年内完全覆盖研发和迁移成本，并持续创造价值。

5 创新点

本文在继承 DADI 成熟的非 FUSE、块级按需加载架构基础上，针对其原有去重机制的局限性，设计并实现了 CORA 方案。其核心创新点如下：

1. **将块级去重从“固定大小”升级为“内容感知”**：我们首次将内容定义分块（Content-Defined Chunking）算法深度集成到 DADI 的块级处理流程中，替代了其原有的固定大小分块策略。这一创新继承了 DADI **块级去重**的思路，但通过从数据内容本身出发识别块边界，有效解决了“边界漂移”问题，将跨层、跨镜像的重复数据识别率提升到新的高度。实验证明，该方法在真实混合数据场景下实现了 **72.36% 的平均空间节省率**。
2. **为非 FUSE 架构设计了高效的变长块随机访问机制**：在 DADI 的 **非 FUSE** 模型下，所有 I/O 都以块的形式在用户态处理。为适应内容分块产生的变长块，我们设计了一套全新的块索引与读取机制。
 - **索引层面**：设计了支持变长块的二级索引结构，并集成了 `raw_block_sizes` 等元数据，实现了逻辑块地址到内容块哈希的高效映射。
 - **读取层面**：重构了 `BlockReader`，通过批量读取、索引追踪和边界保护机制，在变长块场景下实现了 **8.79 倍**的平均解压缩加速，并保证了 **100%** 的数据完整性。
3. **构建了与云原生生态无缝集成的非侵入式加速方案**：我们继承并扩展了 DADI 的 `snapshotter` 插件模式，将 CORA 的全部能力封装为一个独立的 `containerd snapshotter`。这使得 CORA 能够以**完全非侵入**的方式被容器引擎管理和调用，无需修改容器引擎或内核。该方案不仅保留了 DADI 对云原生生态的良好兼容性，还成功完成了向 OpenCloudOS 系统的迁移，证明了其通用性与可移植性。
4. **建立了面向块级按需加载的标准化性能评估框架**：为科学量化内容分块带来的性能收益，我们设计了一套包含三类数据集（随机、模式、混合）、多种文件大小、多次重复实验的基准测试体系。该框架不仅验证了 CORA 的性能，也可作为未来研究块级按需加载技术的标准测试基准。

这些创新共同将 DADI 的**非 FUSE 块级按需加载**架构优势与内容定义分块的**极致去重能力**相结合，实现了“高效去重 + 快速启动 + 数据完整性 + 云原生兼容”的综合目标。

6 应用价值

本方案可广泛应用于：

- **Serverless 平台**：显著减少冷启动延迟，实验显示容器启动时间减少约 50%，满足 Serverless 场景对毫秒级响应的需求；
- **大规模弹性伸缩**：快速拉起海量容器实例，得益于 72% 的空间节省和 40% 的网络流量减少，可大幅降低集群扩容时的资源消耗；

- **混合云环境**：保持跨平台镜像一致性，内容定义分块确保不同环境下的镜像具有相同的哈希指纹，简化镜像分发和同步；
- **边缘计算场景**：按需加载降低带宽占用，在带宽受限的边缘节点上，40% 的流量减少可显著提升镜像分发效率；
- **CI/CD 流水线**：加速构建和测试流程，快速的镜像加载和去重能力可缩短流水线执行时间；
- **容器镜像仓库**：优化存储成本，对于大型镜像仓库，35% 的存储空间节省可带来显著的成本收益（以 1PB 存储为例，年度可节省约 100 万元人民币）。

其在云原生基础设施层面具有明显的产业落地价值，特别适合大规模容器集群和多租户云平台场景。

7 未来工作与展望

7.1 短期优化方向

1. 真实容器镜像测试：

- 扩展到更多真实镜像：Ubuntu、Alpine、Node.js、Python、Java 应用镜像等；
- 测试多层镜像的跨层去重效果；
- 验证镜像更新场景下的增量分块效率。

2. 预取策略优化：

- 基于机器学习的访问模式预测；
- 自适应预取窗口大小调整；
- 热点数据识别与缓存优化。

3. 分块参数调优：

- 针对不同应用类型（数据库、Web 服务、大数据应用）的最优参数配置；
- 动态调整 CDC 参数（最小块、平均块、最大块大小）；
- 探索更高效的哈希算法（如 xxHash、BLAKE3）。

7.2 中期研究方向

1. 跨镜像去重：

- 构建全局块索引，实现镜像仓库级别的去重；
- 探索分层哈希索引结构，优化查询效率；
- 设计增量更新机制，减少索引维护开销。

2. 多级缓存架构：

- L1 缓存（内存）：热点块快速访问；
- L2 缓存（本地 SSD）：常用块持久化；
- L3 缓存（远程对象存储）：全量数据备份；
- 设计智能的缓存替换策略（如 ARC、2Q）。

3. 安全性增强：

- 集成加密分块（Convergent Encryption）；
- 防止哈希碰撞攻击；
- 支持镜像签名与验证。

7.3 长期研究方向

1. 智能化镜像管理：

- 基于深度学习的镜像访问模式预测；
- 自动化镜像优化建议（如合并冗余层、调整分块策略）；
- 镜像质量评估与异常检测。

2. 分布式去重系统：

- 跨数据中心的全局去重；
- 分布式哈希表（DHT）索引；
- 一致性哈希与数据迁移策略。

3. 硬件加速：

- GPU 加速哈希计算；
- FPGA 加速分块与压缩；

- 智能网卡 (SmartNIC) 卸载块传输。

4. 标准化与生态建设:

- 推动 OCI 标准扩展, 纳入内容定义分块规范;
- 开源社区合作, 集成到 containerd、CRI-O 等主流项目;
- 构建性能基准测试套件, 形成行业标准。

7.4 潜在研究问题

- **理论问题:** 最优分块算法的数学模型, 去重率与访问效率的理论上界分析。
- **工程问题:** 超大规模镜像仓库 (PB 级) 的索引构建与维护, 高并发场景下的锁竞争优化。
- **应用问题:** 不同行业 (金融、医疗、制造业) 容器镜像的特征分析与定制化优化。

8 总结

本文针对容器镜像按需加载系统 (DADI) 中固定分块去重效率不足的问题, 提出了基于内容定义分块的优化方案。通过设计变长块索引结构、优化 BlockReader 批量读取机制、构建严格的性能测试框架, 成功实现了高效、可靠的块级镜像管理系统。

主要贡献:

- 在真实混合数据场景下实现了 **72.36% 的平均空间节省率**;
- 实现了 **8.79 倍的平均解压缩加速比**;
- 通过 90 组严格测试验证了 **100% 的数据完整性**;
- 容器启动时间减少约 **50%**;
- 网络流量减少 **40%**, 存储空间节省 **35%**。

参考文献

- [1] Wen Xia, et al. "FastCDC: a Fast and Efficient Content-Defined Chunking Approach for Data Deduplication." USENIX ATC, 2016.
- [2] DADI: Data Accelerator for Disaggregated Infrastructure. Alibaba Cloud, 2020.
- [3] OverlayBD: A Block Device for Container Image. Alibaba Cloud, <https://github.com/containerd/overlaybd>
- [4] OpenCloudOS Community, <https://opencloudos.org>
- [5] Containerd: An industry-standard container runtime. CNCF, <https://containerd.io>
- [6] Open Container Initiative (OCI) Specifications. <https://opencontainers.org>
- [7] Network Block Device (NBD). Linux Kernel Documentation.
- [8] LZ4: Extremely Fast Compression. Yann Collet, <https://lz4.github.io/lz4/>
- [9] Dutch T. Meyer and William J. Bolosky. "A study of practical deduplication." ACM Transactions on Storage (TOS), 2011.

A 实验数据详情

A.1 实验一：容器冷启动性能测试详情

本节实验的原始数据存储于 `cold_start_results.csv`，包含了 OverlayFS 与 OverlayBD 在不同启动场景下的详细耗时记录。测试脚本 `benchmark_cold_start.sh` 用于执行完整的冷启动和热启动测试流程。

```
#!/usr/bin/env bash
set -euo pipefail
# ok!S
IMAGE_BASE=xfusion5:5000/redis:7.2.3
IMAGE_OBD=xfusion5:5000/redis:7.2.3_obd
CTR_BIN=/home/wfy/DADI_OverlayBD_demo/Accelerated_Container_Image/bin/ctr
HOSTS_DIR=/etc/containerd/certs.d
RPULL_HTTPS_PROXY=10.26.42.239:7890
# System-level benchmark: overlaybd remote snapshotter vs standard snapshotter
# Requires containerd + nerdctl (preferred) or ctr. You must provide image refs.
# Metrics: pull time, trivial run time (echo), first I/O time reading a file inside container.

# Config (overridable via env)
SNAP_BASELINE=${SNAP_BASELINE:-overlayfs}
SNAP_OBD=${SNAP_OBD:-overlaybd}
N_RUNS=${N_RUNS:-5}
PROBE_SIZE=${PROBE_SIZE:-10485760} # bytes to read for first-I/O (10MB)
PROBE_PATHS_DEFAULT="/usr/bin/python3_/bin/bash_/lib/x86_64-linux-gnu/libc.so.6_/bin/sh_/usr/bin/busybox"
PROBE_PATHS=${PROBE_PATHS:-$PROBE_PATHS_DEFAULT}
DROP_CACHES=${DROP_CACHES:-1} # 1 to drop host caches (requires sudo)
# Network mode for container runs (to avoid CNI dependency). host|none|default
# host: --net=host (nerdctl) / --net-host (ctr); none: --net=none (nerdctl); default: no extra flag
RUN_NET_MODE=${RUN_NET_MODE:-host}
# Optional: override commands (avoid shell inside image)
```

```
# e.g. RUN_CMD="/usr/local/bin/redis-server --version"
#      IO_CMD="dd if=/usr/local/bin/redis-server of=/dev/null bs=1M
#            count=10"
RUN_CMD=${RUN_CMD:-}
IO_CMD=${IO_CMD:-}

# Concurrency testing config
# CONCURRENCY: number of containers to start concurrently for cold start
#              test
# HOLD_SECONDS: how long each container should stay alive (to observe
#              memory)
CONCURRENCY=${CONCURRENCY:-1000}
HOLD_SECONDS=${HOLD_SECONDS:-30}

# Images (must be provided)
IMAGE_BASE=${IMAGE_BASE:-}
IMAGE_OBD=${IMAGE_OBD:-}

# overlaybd pull via custom ctr (user-provided)
# You can override these to match your environment.
# Example from user:
#   CTR_BIN=/home/wfy/DADI_OverlayBD_demo/Accelerated_Container_Image/
#   bin/ctr \
#   HOSTS_DIR=/etc/containerd/certs.d \
#   RPULL_HTTPS_PROXY=10.26.42.239:7890
CTR_BIN=${CTR_BIN:-/home/wfy/DADI_OverlayBD_demo/
Accelerated_Container_Image/bin/ctr}
HOSTS_DIR=${HOSTS_DIR:-/etc/containerd/certs.d}
RPULL_HTTPS_PROXY=${RPULL_HTTPS_PROXY:-}

# Tool detection
HAVE_NERDCTL=0
HAVE_CTR=0
if command -v nerdctl >/dev/null 2>&1; then HAVE_NERDCTL=1; fi
if command -v ctr >/dev/null 2>&1; then HAVE_CTR=1; fi
if [[ $HAVE_NERDCTL -eq 0 && $HAVE_CTR -eq 0 ]]; then
    echo "[ERROR] Need nerdctl or ctr in PATH" >&2
    exit 1
fi
```



```
fi

OUT_ROOT=$(cd "$(dirname "$0")" && pwd)/out
STAMP=$(date +%Y%m%d-%H%M%S)
OUT_DIR="$OUT_ROOT/$STAMP"
mkdir -p "$OUT_DIR"
RESULT_CSV="$OUT_DIR/sys_results.csv"
echo "mode,snapshotter,image,pull_ms,run_echo_ms,first_io_ms,bytes_rx,
      bytes_tx,notes" > "$RESULT_CSV"

# Extra result files for extended tests
COLD_FULL_CSV="$OUT_DIR/cold_full.csv"
echo "mode,snapshotter,image,cold_full_ms,pull_ms,run_ms,bytes_rx,
      bytes_tx,notes" > "$COLD_FULL_CSV"
CONCUR_CSV="$OUT_DIR/concurrent_results.csv"
echo "mode,snapshotter,image,concurrency,avg_start_ms,min_start_ms,
      max_start_ms,mem_increase_bytes,per_container_bytes,notes" > "
      $CONCUR_CSV"
HOT_CSV="$OUT_DIR/hot_results.csv"
echo "mode,snapshotter,image,run_ms,notes" > "$HOT_CSV"

echo "#_Output_dir:_$OUT_DIR" >&2

require_images() {
    if [[ -z "${IMAGE_BASE}" || -z "${IMAGE_OBD}" ]]; then
        cat >&2 <<EOF
[ERROR] Please set IMAGE_BASE and IMAGE_OBD, e.g.:
    IMAGE_BASE=docker.io/acerak01/overlaybd-image:redis-8.2.1 \
    IMAGE_OBD=docker.io/acerak01/overlaybd-image:latest_obd_new \
    CTR_BIN=/home/wfy/DADI_OverlayBD_demo/Accelerated_Container_Image/bin/
        ctr \
    HOSTS_DIR=/etc/containerd/certs.d \
    RPULL_HTTPS_PROXY=10.26.42.239:7890 \
    bash experiments/sys_bench_overlaybd.sh
EOF
        exit 1
    fi
}
```

```
# Return total RX and TX bytes across all interfaces from /proc/net/dev
net_bytes_total() {
    awk 'NR>2{rx+=$2; tx+=$10} END{printf "%ld %ld\n", rx, tx}' /proc/net/
    dev
}

# Drop host page cache if DROP_CACHES==1 (sudo required)
drop_caches() {
    if [[ "$DROP_CACHES" == "1" ]]; then
        echo "[INFO] Dropping host caches (sudo required)" >&2
        sudo sh -c 'sync; echo 3 > /proc/sys/vm/drop_caches' || true
    fi
}

# Temporarily override DROP_CACHES within a scope
push_drop_caches() {
    __DC_SAVED=${DROP_CACHES:-0}
    DROP_CACHES="1"
}

pop_drop_caches() {
    DROP_CACHES=${__DC_SAVED:-0}
    unset __DC_SAVED
}

# Return host memory used bytes (MemTotal - MemAvailable)
mem_used_bytes() {
    awk '/MemTotal:/{t=$2}/MemAvailable:/{a=$2} END{printf "%ld\n", (t-a)
        *1024}' /proc/meminfo
}

# Pull the image using the selected snapshotter
# - overlaybd: use user-provided ctr rpull with optional proxy and hosts
#   -dir
# - baseline: use nerdctl/ctr standard pull
pull_with() {
    local snap="1"; shift
    local image="1"; shift
}
```

```

local start end dur notes=""

# Remove existing image to force download
if [[ "$snap" == "overlaybd" ]]; then
    # Use user-specified ctr rpull flow for overlaybd images
    if [[ -x "$CTR_BIN" ]]; then
        "$CTR_BIN" i rm "$image" >/dev/null 2>&1 || true
    else
        echo "[WARN]␣CTR_BIN␣not␣executable␣at␣$CTR_BIN;␣trying␣'ctr'␣in␣
            PATH" >&2
        if command -v ctr >/dev/null 2>&1; then ctr images rm "$image" >/
            dev/null 2>&1 || true; fi
    fi
    # Optional: prune nerdctl images to ensure clean state (best-effort)
    if command -v nerdctl >/dev/null 2>&1; then
        nerdctl image prune --force --all >/dev/null 2>&1 || true
    fi
else
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
        nerdctl --snapshotter "$snap" image rm -f "$image" >/dev/null 2>&1
            || true
    else
        ctr images rm "$image" >/dev/null 2>&1 || true
    fi
fi

drop_caches
local rx0 tx0 rx1 tx1
read -r rx0 tx0 < <(net_bytes_total)
start=$(date +%s%N)
if [[ "$snap" == "overlaybd" ]]; then
    # Perform pull using "ctr rpull --hosts-dir" with optional
    https_proxy
    if [[ -x "$CTR_BIN" ]]; then
        if [[ -n "$RPULL_HTTPS_PROXY" ]]; then
            if ! https_proxy="$RPULL_HTTPS_PROXY" "$CTR_BIN" rpull --hosts-
                dir "$HOSTS_DIR" "$image" 1>/dev/null; then
                notes="pull_failed"; dur=-1; rx1=$rx0; tx1=$tx0; echo "$dur␣

```

```
        $rx1_$tx1_$notes"; return
    fi
else
    if ! "$CTR_BIN" rpull --hosts-dir "$HOSTS_DIR" "$image" 1>/dev/
        null; then
        notes="pull_failed"; dur=-1; rx1=$rx0; tx1=$tx0; echo "$dur_
            $rx1_$tx1_$notes"; return
        fi
    fi
else
    # Fallback to standard ctr pull if custom ctr not available
    if ! ctr images pull --snapshotter "$snap" "$image" 1>/dev/null;
        then
        notes="pull_failed"; dur=-1; rx1=$rx0; tx1=$tx0; echo "$dur_$rx1
            _$tx1_$notes"; return
        fi
    fi
else
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
        if ! nerdctl --snapshotter "$snap" image pull "$image" 1>/dev/null
            ; then
            notes="pull_failed"; dur=-1; rx1=$rx0; tx1=$tx0; echo "$dur_$rx1
                _$tx1_$notes"; return
            fi
        else
            if ! ctr images pull --snapshotter "$snap" "$image" 1>/dev/null;
                then
                notes="pull_failed"; dur=-1; rx1=$rx0; tx1=$tx0; echo "$dur_$rx1
                    _$tx1_$notes"; return
                fi
            fi
        fi
    fi
end=$(date +%s%N)
read -r rx1 tx1 < <(net_bytes_total)
dur=$(( (end-start)/1000000 ))
echo "$dur_$rx1_$tx1_$notes"
}
```

```

# Run a trivial command inside container to measure container startup
  overhead
# - Uses RUN_CMD if provided; otherwise runs 'echo ready'
run_echo_with() {
  local snap="$1"; shift
  local image="$1"; shift
  local start end dur notes=""
  drop_caches
  start=$(date +%s%N)
  if [[ $HAVE_NERDCTL -eq 1 ]]; then
    local net_args=""
    case "$RUN_NET_MODE" in
      host) net_args="--net=host";;
      none) net_args="--net=none";;
      *)    net_args="";;
    esac
    if [[ -n "$RUN_CMD" ]]; then
      read -r -a __rcmd <<< "$RUN_CMD"
      nerdctl --snapshotter "$snap" run $net_args --rm "$image" "${
        __rcmd[@]}" >/dev/null 2>&1 || notes="run_failed"
    else
      nerdctl --snapshotter "$snap" run $net_args --rm "$image" sh -c '
        echo_ready' >/dev/null 2>&1 || notes="run_failed"
    fi
  else
    local net_args=""
    case "$RUN_NET_MODE" in
      host) net_args="--net=host";;
      none) net_args="";; # ctr lacks a simple "none" without CNI; use
        default if requested
      *)    net_args="";;
    esac
    if [[ -n "$RUN_CMD" ]]; then
      read -r -a __rcmd <<< "$RUN_CMD"
      ctr run --rm --snapshotter "$snap" $net_args "$image" test-echo "${
        __rcmd[@]}" >/dev/null 2>&1 || notes="run_failed"
    else
      ctr run --rm --snapshotter "$snap" $net_args "$image" test-echo sh

```

```

        -c 'echo_ready' >/dev/null 2>&1 || notes="run_failed"

    fi
fi
end=$(date +%s%N)
dur=$(( (end-start)/1000000 ))
echo "$dur_$notes"
}

# Measure first-read latency by reading PROBE_SIZE bytes from first
# available path in PROBE_PATHS
run_first_io_with() {
    local snap="$1"; shift
    local image="$1"; shift
    local start end dur notes=""
    local ps="$PROBE_SIZE"
    local paths=( $PROBE_PATHS )
    local probe_script='set -e; for p in "${paths[*]}"; do if [ -r "$p" ];
    then echo "Using $p" 1>&2; if command -v dd >/dev/null 2>&1;
    then dd if="$p" of=/dev/null bs=1M count=$((ps/1048576)) status=
    none || head -c "$ps" "$p" >/dev/null; else head -c "$ps" "$p"
    >/dev/null; fi; exit 0; fi; done; echo "No probe file found" 1>&2;
    exit 0'

    drop_caches
    start=$(date +%s%N)
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
        local net_args=""
        case "$RUN_NET_MODE" in
            host) net_args="--net=host";;
            none) net_args="--net=none";;
            *)     net_args="";;
        esac
        if [[ -n "$IO_CMD" ]]; then
            read -r -a __icmd <<< "$IO_CMD"
            nerdctl --snapshotter "$snap" run $net_args --rm "$image" "${
            __icmd[@]}" >/dev/null 2>&1 || notes="io_failed"
        else
            nerdctl --snapshotter "$snap" run $net_args --rm "$image" sh -c "
            $probe_script" >/dev/null 2>&1 || notes="io_failed"
        fi
    fi
}

```

```

    fi
else
    local net_args=""
    case "$RUN_NET_MODE" in
        host) net_args="--net-host";;
        none) net_args="";;
        *)    net_args="";;
    esac
    if [[ -n "$IO_CMD" ]]; then
        read -r -a __icmd <<< "$IO_CMD"
        ctr run --rm --snapshotter "$snap" $net_args "$image" test-io "${__icmd[@]}" >/dev/null 2>&1 || notes="io_failed"
    else
        ctr run --rm --snapshotter "$snap" $net_args "$image" test-io sh -c "$probe_script" >/dev/null 2>&1 || notes="io_failed"
    fi
fi
end=$(date +%s%N)
dur=$(( (end-start)/1000000 ))
echo "$dur_␣$notes"
}

# Measure "single full cold start" time as pull_ms + run_ms
# Steps:
# 1) Forcefully remove image (already done in pull_with)
# 2) Drop caches
# 3) Pull image (timed)
# 4) Immediately run trivial command (timed)
# Output to $COLD_FULL_CSV
bench_single_cold_full() {
    local mode="$1"; shift
    local snap="$1"; shift
    local image="$1"; shift
    echo "[INFO]_␣($mode)_␣Single_␣full_␣cold_␣start_␣for_␣$image" >&2
    push_drop_caches 1
    read -r pull_ms rx1 tx1 notes1 < <(pull_with "$snap" "$image")
    read -r run_ms notes2 < <(run_echo_with "$snap" "$image")
    pop_drop_caches

```

```

local notes="${notes1}${notes2:+,$notes2}"
local total_ms=-1
if [[ $pull_ms -ge 0 && $run_ms -ge 0 ]]; then
    total_ms=$((pull_ms + run_ms))
fi
echo "$mode,$snap,$image,$total_ms,$pull_ms,$run_ms,$rx1,$tx1,$notes"
    | tee -a "$COLD_FULL_CSV" >/dev/null
}

# Launch many containers concurrently and measure average "start latency
"
# - Prefer nerdctl (supports -d). Attempts ctr --detach if available.
# - Each container runs 'sleep HOLD_SECONDS' to keep it alive for memory
    observation.
# - We measure per-launch CLI return time as start latency proxy.
# - Outputs to $CONCUR_CSV: average/min/max start ms and memory increase
    .
bench_concurrent_cold() {
    local mode="$1"; shift
    local snap="$1"; shift
    local image="$1"; shift
    local N=${CONCURRENCY}
    echo "[INFO]␣($mode)␣Concurrent␣cold␣start:␣N=$N␣image=$image" >&2

    # Ensure image is present to isolate start latency (not distribution)
    # Use a warm pull (no delete) to avoid re-download for each container
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
        nerdctl --snapshotter "$snap" image inspect "$image" >/dev/null 2>&1
        || nerdctl --snapshotter "$snap" image pull "$image" >/dev/null
    else
        ctr images list | grep -q "\b$image\b" || ctr images pull --
            snapshotter "$snap" "$image" >/dev/null
    fi

    # Compose network args
    local net_args=""
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
        case "$RUN_NET_MODE" in

```



```
    host) net_args="--net=host";;
    none) net_args="--net=none";;
    *)    net_args="";;
  esac
else
  case "$RUN_NET_MODE" in
    host) net_args="--net=host";;
    *)    net_args="";;
  esac
fi

# Prepare temp dir for timings
local tdir
tdir=$(mktemp -d)
local mem_before mem_after
mem_before=$(mem_used_bytes)

# Spawn N containers
echo "[INFO] Spawning $N containers (hold ${HOLD_SECONDS}s)" >&2
local i
for i in $(seq 1 "$N"); do
  (
    local t0 t1 dur
    t0=$(date +%s%N)
    if [[ $HAVE_NERDCTL -eq 1 ]]; then
      # nerdctl detached run
      nerdctl --snapshotter "$snap" run $net_args -d --rm "$image" sh
        -c "sleep ${HOLD_SECONDS}" >/dev/null 2>&1 || true
    else
      # Attempt ctr detach if supported; else run in background (best-
        effort)
      if ctr run --help 2>/dev/null | grep -q "--detach"; then
        ctr run --rm --detach --snapshotter "$snap" $net_args "$image"
          "c$i" sh -c "sleep ${HOLD_SECONDS}" >/dev/null 2>&1 || true
      else
        ctr run --rm --snapshotter "$snap" $net_args "$image" "c$i" sh
          -c "sleep ${HOLD_SECONDS}" >/dev/null 2>&1 &
        disown || true
      fi
    fi
  )
done
```

```

        fi
    fi
    t1=$(date +%s%N)
    dur=$(( (t1 - t0)/1000000 ))
    echo "$dur" > "$tdir/$i.ms"
) &
done
wait

# Allow processes to settle, then sample memory
sleep 3
mem_after=$(mem_used_bytes)

# Aggregate timings
local sum=0 min=999999999 max=0 count=0 v
for vfile in "$tdir"/*.ms; do
    read -r v < "$vfile" || v=0
    sum=$((sum + v))
    (( v < min )) && min=$v
    (( v > max )) && max=$v
    count=$((count + 1))
done
local avg=0
if [[ $count -gt 0 ]]; then avg=$((sum / count)); fi
local mem_inc=$((mem_after - mem_before))
local per_ct=0
if [[ $N -gt 0 ]]; then per_ct=$(( mem_inc / N )); fi

echo "$mode,$snap,$image,$N,$avg,$min,$max,$mem_inc,$per_ct," | tee -a
"$CONCUR_CSV" >/dev/null
rm -rf "$tdir"
}

# Measure hot start run time (no cache drops), image must be present
already
bench_hot_start() {
    local mode="$1"; shift
    local snap="$1"; shift

```

```

local image="$1"; shift
echo "[INFO]␣($mode)␣Hot␣start␣runs␣for␣$image␣(N=$N_RUNS)" >&2

# Ensure image present
if [[ $HAVE_NERDCTL -eq 1 ]]; then
    nerdctl --snapshotter "$snap" image inspect "$image" >/dev/null 2>&1
    || nerdctl --snapshotter "$snap" image pull "$image" >/dev/null
else
    ctr images list | grep -q "\b$image\b" || ctr images pull --
        snapshotter "$snap" "$image" >/dev/null
fi

push_drop_caches 0
local i
for i in $(seq 1 "$N_RUNS"); do
    read -r run_ms notes < <(run_echo_with "$snap" "$image")
    echo "$mode,$snap,$image,$run_ms,$notes" | tee -a "$HOT_CSV" >/dev/
        null
done
pop_drop_caches
}

bench_mode() {
    local mode="$1"; shift
    local snap="$1"; shift
    local image="$1"; shift
    for i in $(seq 1 "$N_RUNS"); do
        echo "[INFO]␣($mode)␣Pull␣$image␣[run␣$i/$N_RUNS]" >&2
        read -r pull_ms rx1 tx1 notes < <(pull_with "$snap" "$image") #
            bugging
        echo "[INFO]␣($mode)␣Run␣echo" >&2
        read -r echo_ms notes2 < <(run_echo_with "$snap" "$image")
        echo "[INFO]␣($mode)␣First␣IO" >&2
        read -r io_ms notes3 < <(run_first_io_with "$snap" "$image")
        local combined_notes="${notes}${notes2:+,$notes2}${notes3:+,$notes3}"
            "
        echo "$mode,$snap,$image,$pull_ms,$echo_ms,$io_ms,$((rx1)),$((tx1)),
            $combined_notes" | tee -a "$RESULT_CSV" >/dev/null
    done
}

```

```
done
}

main() {
    require_images
    bench_mode baseline "$SNAP_BASELINE" "$IMAGE_BASE"
    bench_mode overlaybd "$SNAP_OBD" "$IMAGE_OBD"

    # Single full cold start (pull + run)
    bench_single_cold_full baseline "$SNAP_BASELINE" "$IMAGE_BASE"
    bench_single_cold_full overlaybd "$SNAP_OBD" "$IMAGE_OBD"

    Concurrent cold start (N containers launching concurrently)
    bench_concurrent_cold baseline "$SNAP_BASELINE" "$IMAGE_BASE"
    bench_concurrent_cold overlaybd "$SNAP_OBD" "$IMAGE_OBD"
    Hot start comparison (no cache drop)
    bench_hot_start baseline "$SNAP_BASELINE" "$IMAGE_BASE"
    bench_hot_start overlaybd "$SNAP_OBD" "$IMAGE_OBD"

    ln -sf "$OUT_DIR" "$OUT_ROOT/latest"
    echo "Done. Results: $RESULT_CSV" >&2
}

main "$@"
```

A.2 实验二：测试脚本说明

基准测试脚本 `benchmark_fastcdc_advanced.sh` 的主要功能模块包括：

- `generate_test_data()`: 生成三类测试数据（随机、模式、混合）
- `benchmark_test_repeated()`: 执行单次基准测试并记录结果
- `calculate_statistics()`: 使用 Python 计算统计量
- `init_csv()`: 初始化 CSV 文件头
- `append_to_csv()`: 追加测试结果到 CSV 文件

```
#!/bin/bash

# 内容定义分块 vs 固定大小分块性能对比测试
# 支持多组重复测试和CSV输出

set -e

# ===== 配置参数 =====
ZFILE=/home/wfy/DADI_OverlayBD_demo/overlaybd/build/output/overlaybd-
zfile
WORK_DIR=/home/wfy/DADI_OverlayBD_demo/experiments/fastcdc_benchmark
RESULT_FILE="$WORK_DIR/benchmark_results.txt"
CSV_FILE="$WORK_DIR/benchmark_results.csv"
CSV_SUMMARY="$WORK_DIR/benchmark_summary.csv"

# 重复测试次数（建议3-5次以获得稳定结果）
REPEAT_COUNT=5

# 颜色输出
RED='\033[0;31m'
GREEN='\033[0;32m'
YELLOW='\033[1;33m'
BLUE='\033[0;34m'
CYAN='\033[0;36m'
NC='\033[0m' # No Color

# 创建工作目录
mkdir -p "$WORK_DIR"
cd "$WORK_DIR"

# 初始化CSV文件
init_csv() {
    # 详细数据CSV（每次测试的原始数据）
    echo "Timestamp,RunID,TestName,DataType,FileSizeKB,OriginalBytes,
Method,CompressedBytes,CompressionRatio,CompressTime_ms,
DecompressTime_ms,MD5Verified" > "$CSV_FILE"
```

```
# 汇总CSV (包含统计数据: 平均值、标准差)
echo "TestName,DataType,FileSizeKB,OriginalBytes,Method,
    AvgCompressedBytes,StdCompressedBytes,AvgCompressionRatio,
    StdCompressionRatio,AvgCompressTime_ms,StdCompressTime_ms,
    AvgDecompressTime_ms,StdDecompressTime_ms,SpaceSaving%,
    SpeedupCompress,SpeedupDecompress" > "$CSV_SUMMARY"
}

# 将测试结果追加到CSV
append_to_csv() {
    local timestamp=$1
    local run_id=$2
    local test_name=$3
    local data_type=$4
    local file_size_kb=$5
    local original_bytes=$6
    local method=$7
    local compressed_bytes=$8
    local compression_ratio=$9
    local compress_time=${10}
    local decompress_time=${11}
    local md5_verified=${12}

    echo "$timestamp,$run_id,$test_name,$data_type,$file_size_kb,
        $original_bytes,$method,$compressed_bytes,$compression_ratio,
        $compress_time,$decompress_time,$md5_verified" >> "$CSV_FILE"
}

# 计算统计数据 (平均值和标准差)
calculate_statistics() {
    python3 << 'PYTHON_EOF'
import csv
import sys
from collections import defaultdict
import math

# 读取CSV数据
data = defaultdict(lambda: defaultdict(list))
```

```
with open('benchmark_results.csv', 'r') as f:
    reader = csv.DictReader(f)
    for row in reader:
        key = (row['TestName'], row['DataType'], row['FileSizeKB'], row[
            'Method'])
        data[key]['CompressedBytes'].append(float(row['CompressedBytes '
            ]))
        data[key]['CompressionRatio'].append(float(row['CompressionRatio
            ']))
        data[key]['CompressTime_ms'].append(float(row['CompressTime_ms '
            ]))
        data[key]['DecompressTime_ms'].append(float(row['
            DecompressTime_ms']))
        data[key]['OriginalBytes'] = row['OriginalBytes']

# 计算平均值和标准差
def mean(values):
    return sum(values) / len(values) if values else 0

def std_dev(values):
    if len(values) < 2:
        return 0
    avg = mean(values)
    variance = sum((x - avg) ** 2 for x in values) / (len(values) - 1)
    return math.sqrt(variance)

# 按测试分组计算
results = defaultdict(dict)
for key, metrics in data.items():
    test_name, data_type, file_size, method = key
    test_key = (test_name, data_type, file_size)

    results[test_key][method] = {
        'avg_compressed': mean(metrics['CompressedBytes']),
        'std_compressed': std_dev(metrics['CompressedBytes']),
        'avg_ratio': mean(metrics['CompressionRatio']),
        'std_ratio': std_dev(metrics['CompressionRatio']),
```

```
'avg_compress_time': mean(metrics['CompressTime_ms']),
'std_compress_time': std_dev(metrics['CompressTime_ms']),
'avg_decompress_time': mean(metrics['DecompressTime_ms']),
'std_decompress_time': std_dev(metrics['DecompressTime_ms']),
'original_bytes': metrics['OriginalBytes']
}

# 写入汇总CSV
with open('benchmark_summary.csv', 'w', newline='') as f:
    writer = csv.writer(f)
    writer.writerow(['TestName', 'DataType', 'FileSizeKB', '
        OriginalBytes', 'Method',
            'AvgCompressedBytes', 'StdCompressedBytes', '
            AvgCompressionRatio', 'StdCompressionRatio',
            'AvgCompressTime_ms', 'StdCompressTime_ms', '
            AvgDecompressTime_ms', 'StdDecompressTime_ms',
            'SpaceSaving%', 'SpeedupCompress', '
            SpeedupDecompress'])

for test_key in sorted(results.keys()):
    test_name, data_type, file_size = test_key

    if 'Fixed' in results[test_key] and 'FastCDC' in results[
        test_key]:
        fixed = results[test_key]['Fixed']
        fastcdc = results[test_key]['FastCDC']

        # 计算节省和加速比
        space_saving = ((fixed['avg_compressed'] - fastcdc['
            avg_compressed']) /
            fixed['avg_compressed'] * 100) if fixed['
            avg_compressed'] > 0 else 0
        speedup_compress = (fixed['avg_compress_time'] / fastcdc['
            avg_compress_time']) if fastcdc['avg_compress_time'] > 0
            else 0
        speedup_decompress = (fixed['avg_decompress_time'] / fastcdc[
            'avg_decompress_time']) if fastcdc['avg_decompress_time']
            > 0 else 0
```



```
# 写入固定分块数据
writer.writerow([
    test_name, data_type, file_size, fixed['original_bytes'
    ], 'Fixed',
    f"{fixed['avg_compressed']:.2f}", f"{fixed['
    std_compressed']:.2f}",
    f"{fixed['avg_ratio']:.2f}", f"{fixed['std_ratio']:.2f}"
    ,
    f"{fixed['avg_compress_time']:.2f}", f"{fixed['
    std_compress_time']:.2f}",
    f"{fixed['avg_decompress_time']:.2f}", f"{fixed['
    std_decompress_time']:.2f}",
    '-', '-', '-'
])

# 写入FastCDC数据
writer.writerow([
    test_name, data_type, file_size, fastcdc['original_bytes
    '], 'FastCDC',
    f"{fastcdc['avg_compressed']:.2f}", f"{fastcdc['
    std_compressed']:.2f}",
    f"{fastcdc['avg_ratio']:.2f}", f"{fastcdc['std_ratio
    ']:.2f}",
    f"{fastcdc['avg_compress_time']:.2f}", f"{fastcdc['
    std_compress_time']:.2f}",
    f"{fastcdc['avg_decompress_time']:.2f}", f"{fastcdc['
    std_decompress_time']:.2f}",
    f"{space_saving:.2f}", f"{speedup_compress:.2f}", f"{
    speedup_decompress:.2f}"
])

print(" 统计分析报告完成")
PYTHON_EOF
}

# 主测试函数（支持重复测试）
benchmark_test_repeated() {
```

```
local test_name=$1
local file_size=$2
local data_type=$3
local run_id=$4
local file_name="${test_name}_${file_size}_run${run_id}"

# 只在第一次运行时显示标题
if [ $run_id -eq 1 ]; then
    echo -e "${BLUE}>>> 测试: ${test_name} (大小: ${file_size}KB, 类型: ${data_type})${NC}"
    echo -e "${CYAN} 重复 ${REPEAT_COUNT} 次以确保结果稳定性...${NC}"
fi

echo -n " 运行 ${run_id}/${REPEAT_COUNT}: "

#bug in

# 生成测试数据
case $data_type in
    "random")
        dd if=/dev/urandom of="${file_name}.dat" bs=1024 count=
            $file_size 2>/dev/null
        ;;
    "pattern")
        python3 << EOF
with open('${file_name}.dat', 'wb') as f:
    pattern = b"ABCDEFGHJKLMNOPQRSTUVWXYZ0123456789" * 100
    total_size = $file_size * 1024
    while f.tell() < total_size:
        f.write(pattern)
EOF
        ;;
    "mixed")
        python3 << EOF
import os
with open('${file_name}.dat', 'wb') as f:
```

```

pattern = b"REPEATED_BLOCK_" * 64
random_block = os.urandom(1024)
total_size = $file_size * 1024
while f.tell() < total_size:
    f.write(pattern)
    f.write(random_block)
EOF

;;

esac

local original_size=$(stat -c%s "${file_name}.dat")
local orig_md5=$(md5sum "${file_name}.dat" | cut -d'_' -f1)
local timestamp=$(date +"%Y-%m-%d_%H:%M:%S")

# === 测试固定分块 ===
local fixed_start=$(date +%s%N)
$ZFILE "${file_name}.dat" "${file_name}_fixed.zfile" 2>/dev/null
local fixed_compress_time=$(( ($date +%s%N) - $fixed_start) /
    1000000 ))
local fixed_size=$(stat -c%s "${file_name}_fixed.zfile")

local fixed_decomp_start=$(date +%s%N)
$ZFILE -x "${file_name}_fixed.zfile" "${file_name}_fixed_out.dat"
    2>/dev/null
local fixed_decompress_time=$(( ($date +%s%N) - $fixed_decomp_start
    ) / 1000000 ))

local fixed_md5=$(md5sum "${file_name}_fixed_out.dat" | cut -d'_' -
    f1)
local fixed_verified="true"
if [ "$fixed_md5" != "$orig_md5" ]; then
    fixed_verified="false"
    echo -e "${RED}固定分块MD5失败!${NC}"
    return 1
fi

local fixed_ratio=$(echo "scale=4;_${fixed_size}_*_${100}_/_
    $original_size" | bc)

```

```
# 保存固定分块数据到CSV
append_to_csv "$timestamp" "$run_id" "$test_name" "$data_type" "
    $file_size" \
    "$original_size" "Fixed" "$fixed_size" "$fixed_ratio" \
    "$fixed_compress_time" "$fixed_decompress_time" "$fixed_verified
    "

# === 测试FastCDC ===
local fastcdc_start=$(date +%s%N)
$ZFILE --fastcdc "${file_name}.dat" "${file_name}_fastcdc.zfile" 2>/
dev/null
local fastcdc_compress_time=$(( (date +%s%N) - $fastcdc_start) /
    1000000 ))
local fastcdc_size=$(stat -c%s "${file_name}_fastcdc.zfile")

local fastcdc_decomp_start=$(date +%s%N)
$ZFILE -x "${file_name}_fastcdc.zfile" "${file_name}_fastcdc_out.dat
    " 2>/dev/null
local fastcdc_decompress_time=$(( (date +%s%N) -
    $fastcdc_decomp_start) / 1000000 ))

local fastcdc_md5=$(md5sum "${file_name}_fastcdc_out.dat" | cut -d' '
    -f1)
local fastcdc_verified="true"
if [ "$fastcdc_md5" != "$orig_md5" ]; then
    fastcdc_verified="false"
    echo -e "${RED}FastCDC MD5 失败!${NC}"
    return 1
fi

local fastcdc_ratio=$(echo "scale=4; ${fastcdc_size}*100/
    $original_size" | bc)

# 保存FastCDC数据到CSV
append_to_csv "$timestamp" "$run_id" "$test_name" "$data_type" "
    $file_size" \
    "$original_size" "FastCDC" "$fastcdc_size" "$fastcdc_ratio" \
```

```
        "$fastcdc_compress_time" "$fastcdc_decompress_time" "
        $fastcdc_verified"

echo -e "${GREEN} ${NC}␣(固定:␣${fixed_compress_time}ms/${
    fixed_decompress_time}ms,␣FastCDC:␣${fastcdc_compress_time}ms/${
    fastcdc_decompress_time}ms)"

# 清理临时文件
rm -f "${file_name}.dat" "${file_name}"_*.zfile,dat}
}

# 运行测试套件
run_test_suite() {
    local test_name=$1
    local file_size=$2
    local data_type=$3

    for run in $(seq 1 $REPEAT_COUNT); do
        benchmark_test_repeated "$test_name" "$file_size" "$data_type" "
            $run"
    done
    echo ""
}

# ===== 主程序 =====

echo "===== " | tee "$RESULT_FILE"
echo "FastCDC␣vs␣固定分块性能对比测试" | tee -a "$RESULT_FILE"
echo "===== " | tee -a "$RESULT_FILE"
echo "测试时间:␣$(date)" | tee -a "$RESULT_FILE"
echo "重复次数:␣$REPEAT_COUNT" | tee -a "$RESULT_FILE"
echo "CSV输出:␣$CSV_FILE" | tee -a "$RESULT_FILE"
echo "统计摘要:␣$CSV_SUMMARY" | tee -a "$RESULT_FILE"
echo "" | tee -a "$RESULT_FILE"

# 初始化CSV文件
init_csv
```

```
echo -e "${YELLOW}===== ${NC}"
echo -e "${YELLOW}测试_1: 随机数据（不可压缩） ${NC}"
echo -e "${YELLOW}===== ${NC}"
run_test_suite "random_100k" 100 "random"
run_test_suite "random_1m" 1024 "random"
run_test_suite "random_10m" 10240 "random"

echo -e "${YELLOW}===== ${NC}"
echo -e "${YELLOW}测试_2: 重复模式数据（高可压缩） ${NC}"
echo -e "${YELLOW}===== ${NC}"
run_test_suite "pattern_100k" 100 "pattern"
run_test_suite "pattern_1m" 1024 "pattern"
run_test_suite "pattern_10m" 10240 "pattern"

echo -e "${YELLOW}===== ${NC}"
echo -e "${YELLOW}测试_3: 混合数据（模拟真实场景） ${NC}"
echo -e "${YELLOW}===== ${NC}"
run_test_suite "mixed_100k" 100 "mixed"
run_test_suite "mixed_1m" 1024 "mixed"
run_test_suite "mixed_10m" 10240 "mixed"

echo -e "${YELLOW}===== ${NC}"
echo -e "${YELLOW}正在计算统计数据... ${NC}"
echo -e "${YELLOW}===== ${NC}"
calculate_statistics

echo ""
echo -e "${GREEN}===== ${NC}"
echo -e "${GREEN}测试完成! ${NC}"
echo -e "${GREEN}===== ${NC}"
echo ""
echo -e "详细数据: ${CYAN}$CSV_FILE ${NC}"
echo -e "统计摘要: ${CYAN}$CSV_SUMMARY ${NC}"
echo ""
echo -e "${YELLOW}统计摘要预览: ${NC}"
echo "-----"

# 显示汇总表格
```

```
python3 << 'PYTHON_EOF'
import csv

with open('benchmark_summary.csv', 'r') as f:
    reader = csv.DictReader(f)
    print(f"{'测试名称':<20}{'方法':<10}{'平均压缩率':<12}{'空间节省':<12}{'解压加速':<12}")
    print("-" * 80)

    for row in reader:
        if row['Method'] == 'FastCDC':
            print(f"{row['TestName']:<20}{row['Method']:<10}"
                  f"{row['AvgCompressionRatio']:>6}%±{row['StdCompressionRatio']:>5}%"
                  f"{row['SpaceSaving%']:>8}%"
                  f"{row['SpeedupDecompress']:>8}x")

PYTHON_EOF

echo ""
echo -e "${CYAN}提示: 使用 Excel/LibreOffice 打开 CSV 文件进行详细分析${NC}"
echo ""

# 返回原始目录
cd - >/dev/null
```

A.3 实验二：完整基准测试结果

完整的 90 组测试原始数据存储于 `benchmark_results.csv`, 统计汇总数据存储于 `benchmark_summary.csv`。所有数据均可通过以下命令查看：

```
# 查看冷启动测试原始数据
cat cold_start_results.csv

# 执行冷启动基准测试
# (需要配置好 containerd 与 snapshotter)
sudo ./benchmark_cold_start.sh
```

A.4 实验二：数据完整性验证

所有测试文件的 MD5 校验结果均为 MD5Verified=true, 证明解压缩后的数据与原始数据完全一致。校验命令示例：

```
# 原始文件 MD5
md5sum test_data.bin

# 压缩后解压缩文件 MD5
./overlaybd-zfile compress --fastcdc test_data.bin compressed.zfile
./overlaybd-zfile decompress compressed.zfile recovered.bin
md5sum recovered.bin

# 比对结果
diff test_data.bin recovered.bin
```