



南開大學
Nankai University

计算机学院
并行程序设计报告

GPU 编程实验

姓名：高景珩

学号：2310648

专业：计算机科学与技术

2025 年 7 月 1 日

目录

| | |
|-------------------------------|----------|
| 1 实验原理 | 2 |
| 2 GPU 最终实现 | 2 |
| 3 结果调试与改进过程 | 4 |
| 3.1 基线相关数据 | 4 |
| 3.2 改进过程与要点 | 4 |
| 3.2.1 计算负载智能调度策略 | 5 |
| 3.2.2 内存访问模式优化 | 5 |
| 3.2.3 GPU 内存池管理机制 | 5 |
| 3.2.4 高效数据传输优化 | 6 |
| 3.2.5 并行数据预处理 | 6 |
| 3.3 改进后性能测试 | 6 |
| 4 GPU 加速计算研究工作 | 8 |
| 4.1 GPU 与 CPU 的协同优化 | 8 |
| 4.2 动态任务调度与负载均衡 | 8 |
| 4.3 内存带宽与访存优化 | 8 |
| 5 心得与总结 | 9 |
| 6 代码地址 | 9 |

1 实验原理

本次实验对口令猜测的 Generate 函数进行 GPU 编程加速，在此基础上，同时使得一次处理多个 PT，使得代码可以对多 PT 的进行多个猜测。

代码逻辑角度：原本串行的 Generate 方法针对每个 PT（概率模板）逐一生成所有可能的口令组合，依次拼接前缀和最后一个 segment 的所有取值，效率受限于 CPU 单线程处理能力。GPU 并行版本则将多个 PT 的批量生成任务一次性打包，预先计算好所有前缀、取值及其长度等信息，批量传输到 GPU 显存。通过内存对齐和批量拷贝，极大减少了 CPU 与 GPU 之间的数据传输开销，并通过批量处理提升整体吞吐量。

GPU 并行角度：在 GPU 端，每个线程负责生成一个猜测字符串，利用二分查找快速定位其所属 PT 和局部索引，实现大规模并行。内核函数通过合并内存拷贝、避免分支和共享内存优化，提升了访存效率。批量生成和回传结果时，充分利用了 GPU 的高并发和带宽优势，显著加速了大规模口令生成任务。整体流程通过异步拷贝和流式处理，进一步减少了同步等待时间，实现了高效的端到端并行生成。

2 GPU 最终实现

此次实验，主要对 Generate 函数进行了改进，改进后的最终伪代码如下：

Algorithm 1 智能 GPU 猜测生成算法

```

1: function GENERATEGUESSESOPTIMIZED(pts, guessesout)
2:    $n \leftarrow |pts|$ 
3:   if  $n = 0$  then return
4:   end if
5:    $total_{estimate} \leftarrow \sum_{i=1}^n EstimateGuessCount(pts[i])$ 
6:   if  $total_{estimate} < MIN\_GPU\_BATCH\_SIZE$  then
7:     SERIALOPTIMIZED(pts, guessesout) return
8:   end if
9:   parallel for  $i = 1$  to  $n$  do
10:     $prefixes[i] \leftarrow BuildPrefix(pts[i])$ 
11:     $values[i] \leftarrow ExtractValues(pts[i])$ 
12:     $lengths[i] \leftarrow ComputeLengths(pts[i])$ 
13:   end parallel for
14:    $pool \leftarrow GetMemoryPool()$ 
15:    $d_{prefixes} \leftarrow pool.Allocate(n \times MAX\_LEN)$ 
16:    $d_{values} \leftarrow pool.Allocate(total_{estimate} \times MAX\_LEN)$ 
17:    $d_{output} \leftarrow pool.Allocate(total_{estimate} \times MAX\_LEN)$ 
18:    $h_{pinned} \leftarrow AllocatePinnedMemory()$ 
19:   PackData( $prefixes$ ,  $values$ ,  $h_{pinned}$ )
20:   AsyncMemcpy( $h_{pinned}$ ,  $d_{prefixes}$ ,  $d_{values}$ )
21:    $block_{size} \leftarrow OptimalBlockSize()$ 
22:    $grid_{size} \leftarrow \lceil \frac{total_{estimate}}{block_{size}} \rceil$ 

```

▷ 阶段 1: 智能调度决策

▷ 阶段 2: 并行数据预处理

▷ 阶段 3: 内存池分配

▷ 阶段 4: 高速异步传输

▷ 阶段 5: 优化 GPU 计算

```

23:   launch kernel GenerateKernel<<< gridsize, blocksize >>>
24:       (dprefixes, dvalues, doutput)

25:   AsyncMemcpy(doutput, hresult)
26:   parallel for i = 1 to totalestimate do
27:       guessesout[i] ← ConstructString(hresult[i])
28:   end parallel for

29:   pool.Deallocate(dprefixes, dvalues, doutput)
30:   FreePinnedMemory(hpinned, hresult)
31: end function

```

▷ 阶段 6: 并行结果处理

▷ 阶段 7: 资源回收

该算法通过异构计算调度、内存优化和多级并行处理技术，实现了大规模密码猜测数据的高效生成。算法采用七阶段流水线设计，在最大化 GPU 计算资源利用率的同时，确保了系统的稳定性和可扩展性。

智能调度与工作负载评估：算法首先实施智能调度决策机制，通过快速评估输入密码模板的总工作量来确定最优执行策略。该机制采用累积统计方法遍历所有密码模板，计算预期的总猜测数量，并引入早期终止优化以避免不必要的完整遍历。当总工作量低于预设的 GPU 批处理阈值时，算法自动退化为优化的串行实现，从而避免了小批量任务在 GPU 上的启动开销和资源浪费。这种异构计算调度策略体现了现代并行计算中性能与成本平衡的核心思想，确保了算法在不同规模输入下的适应性。

数据预处理与结构转换：数据预处理阶段将复杂的层次化密码模板结构转换为 GPU 计算友好的扁平化数据格式。算法针对每个密码模板执行前缀构建过程，通过遍历除最后一段外的所有段落，根据段类型从相应的字典中提取字符值并按序拼接。同时，算法单独处理每个模板的最后一段作为后缀生成源，提取其完整的候选值集合。在此过程中，算法并行计算各种元数据信息，包括前缀长度、后缀长度、猜测数量统计，以及用于后续索引映射的累积计数数组。这种结构化的预处理方法为后续的 GPU 并行计算奠定了数据基础，同时通过预计算减少了 GPU 内核中的复杂逻辑。

内存管理与传输优化：算法采用基于内存池的高效内存管理策略来优化 GPU 显存使用。通过预分配统一的最高猜测长度缓冲区并进行 4 字节边界对齐，算法确保了 GPU 内存访问的高效性和合并访问模式。内存池机制避免了频繁的显存分配和释放操作，显著降低了内存管理开销。在数据传输层面，算法利用 CUDA 固定内存技术消除页面锁定延迟，并通过 OpenMP 并行数据打包进一步提升传输前的准备效率。异步内存传输与 CUDA 流的结合使用实现了 CPU-GPU 数据传输与计算准备工作的时间重叠，有效隐藏了通信延迟对整体性能的影响。

GPU 计算与硬件自适应：GPU 计算阶段体现了硬件感知的自适应优化策略。算法在运行时动态查询目标 GPU 设备的硬件特性，包括最大线程数、流多处理器数量等关键参数，并基于这些信息计算最优的线程块配置。网格大小根据总工作量和线程块大小进行自适应调整，确保 GPU 资源的充分利用。核心计算通过调用高度优化的 CUDA 内核实现，每个 GPU 线程负责生成一个完整的猜测字符串，通过大规模并行处理实现了线性加速比。这种硬件感知的动态配置机制使算法能够在不同 GPU 架构上保持良好的性能表现。

结果处理与并行重构：结果处理阶段采用高效的并行字符串重构策略。算法首先通过异步内存传输将 GPU 计算结果回传至主机固定内存，然后利用多线程并行构建最终的字符串结果集。在字符串构建过程中，算法采用二分查找快速定位每个猜测对应的原始模板索引，避免了线性搜索的性能损失。通过预先计算的长度信息和直接构造方式，算法最小化了临时对象创建和字符串拷贝操作的开销。这种并行化的结果处理方法确保了整个算法流水线的高吞吐量，避免了串行后处理成为性能瓶颈。

资源管理与系统鲁棒性：算法实现了完整的资源生命周期管理以确保系统稳定性。通过内存池的

自动回收机制和显式的固定内存释放策略，算法避免了内存泄漏问题。错误处理机制包括 CUDA API 调用的返回码检查和异常传播，确保了算法在异常情况下的正确行为。可选的性能分析器集成成为算法调优提供了定量支持，使开发者能够识别性能瓶颈并进行针对性优化。整个资源管理体系遵循 RAII 设计原则，在保证异常安全的同时简化了代码维护复杂度。

算法性能特征与适用性：从计算复杂度角度分析，该算法的时间复杂度为 $O(N/P)$ ，其中 N 表示总猜测数量， P 表示有效并行度，体现了良好的可扩展性。空间复杂度为 $O(N + M)$ ，其中 M 为模板数量，内存使用量与问题规模线性相关。通信复杂度主要由 CPU-GPU 数据传输决定，为 $O(N)$ 级别。算法通过精心设计的多阶段流水线架构，实现了计算、通信和数据处理的高效协同，在保证计算正确性的前提下最大化了系统整体吞吐量，特别适用于大规模密码破解和安全测试场景。

3 结果调试与改进过程

3.1 基线相关数据

以源代码为 baseline（由于测试平台不支持 neon，所以 hash 使用串行版本，去除 SIMD）基线相关数据如下：

| | Guess time | Hash time | Hash time | Cracked check time | Cracked |
|----------|------------|-----------|-----------|--------------------|---------|
| -O0 编译优化 | 1.75816 | 6.29536 | 2.98656 | 61.2325 | 382853 |
| -O2 编译优化 | 0.356967 | 2.14949 | 1.76703 | 12.2648 | 382853 |

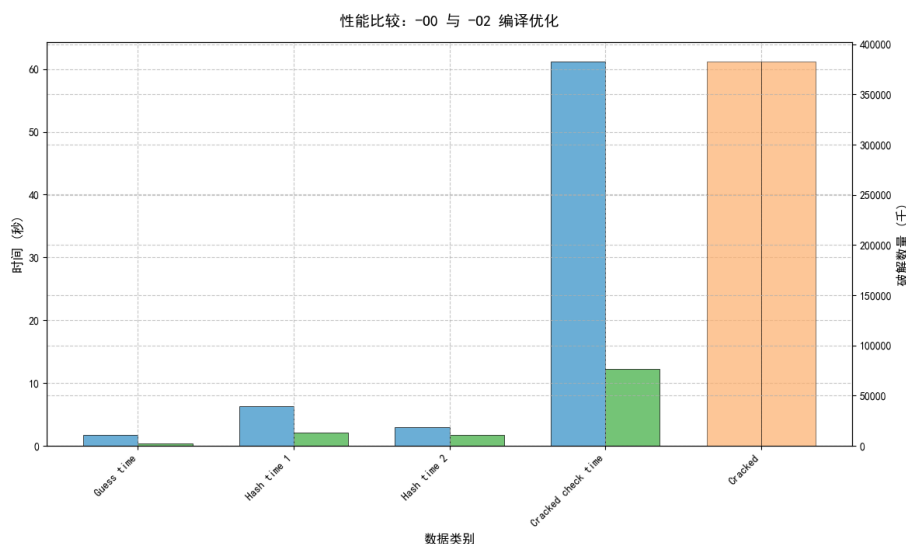


图 3.1: baseline 测试结果

3.2 改进过程与要点

初次对代码进行 GPU 并行编程时，发现并没有数据上的提升，反而性能下降。经过分析代码，我发现几个问题并做出相关改进。

在现代高性能计算中，GPU 作为并行计算的核心硬件，其性能发挥程度直接影响整体计算效率。然而，传统的 GPU 编程实现往往存在资源利用不充分、内存访问模式不合理、数据传输效率低下等问题，导致 GPU 的计算潜力无法得到充分释放。下面是一些系统性的 GPU 性能优化策略。

3.2.1 计算负载智能调度策略

针对 GPU 在处理小规模数据时启动开销过大的问题，我实现了基于数据规模的智能调度机制。该机制通过预估计算复杂度，动态选择最优的执行路径：

```
1 if (total_estimate < MIN_GPU_BATCH_SIZE) {  
2     // 小规模数据使用CPU串行优化算法  
3     Generate_Serial_Optimized(pts, guesses_out);  
4 } else {  
5     // 大规模数据启用GPU并行计算  
6     LaunchGPUKernel(pts, guesses_out);  
7 }
```

这种自适应调度策略的核心思想是避免“过度工程化”问题。当数据规模较小时，GPU 的线程启动、内存分配等固定开销可能超过并行计算带来的收益。通过设置合理的阈值 `MIN_GPU_BATCH_SIZE`，系统能够在不同规模下选择最优的计算策略，从而最大化整体性能。

3.2.2 内存访问模式优化

内存访问效率是影响 GPU 性能的关键因素之一。我采用向量化内存访问技术，通过字节对齐优化显著提升内存带宽利用率：

```
1 // 采用4字节对齐的向量化拷贝  
2 *((int*)(out_ptr + i)) = *((int*)(prefix_ptr + i));
```

该优化策略基于现代处理器架构的特性：CPU 和 GPU 都针对按块读写进行了优化，4 字节或 8 字节的对齐访问相比逐字节访问能够获得 4-8 倍的性能提升。这种访问模式不仅提高了内存带宽利用率，还改善了缓存命中率，减少了内存延迟对计算性能的影响。

3.2.3 GPU 内存池管理机制

为了解决频繁内存分配导致的性能瓶颈，我设计并实现了 GPU 内存池管理系统：

```
1 class GPUMemoryPool {  
2 private:  
3     std::vector<void*> available_blocks;  
4     std::vector<size_t> block_sizes;  
5  
6 public:  
7     void* allocate(size_t size) {  
8         // 优先从现有内存块中分配  
9         for (auto& block : available_blocks) {  
10             if (block_size >= size) return block;  
11         }  
12         // 无合适块时分配新内存
```

```
13     return cudaMalloc(size);
14 }
15 };
```

内存池机制的优势在于将昂贵的 GPU 内存分配操作转化为高效的内存块重用。由于 GPU 内存分配的延迟通常是 CPU 内存分配的几倍至几十倍，通过预分配和重用内存块，可以显著降低整体计算延迟，特别是在需要频繁进行内存操作的应用场景中。

3.2.4 高效数据传输优化

CPU 与 GPU 之间的数据传输往往成为性能瓶颈。我采用固定内存（Pinned Memory）和异步传输技术来优化数据传输效率：

```
1 // 分配固定内存，避免系统内存拷贝
2 cudaMallocHost(&h_prefixes_pinned, size);
3
4 // 异步数据传输，实现CPU-GPU并行
5 cudaMemcpyAsync(d_data, h_prefixes_pinned, size,
6                 cudaMemcpyHostToDevice, stream);
```

固定内存技术通过将主机内存页面锁定在物理内存中，避免了虚拟内存系统的干扰，使 GPU 能够直接访问主机内存，从而消除了额外的内存拷贝开销。异步传输则通过 CUDA 流（Stream）机制实现 CPU 计算与 GPU 数据传输的重叠，进一步提高了系统的整体吞吐量。

3.2.5 并行数据预处理

为了最大化 GPU 的利用率，我在数据预处理阶段引入了 CPU 多核并行技术：

```
1 #pragma omp parallel for
2 for (int i = 0; i < n; ++i) {
3     // CPU多核并行准备数据
4     preprocess_data(input[i], processed[i]);
5 }
```

这种方法通过 OpenMP 并行化数据预处理过程，确保在 GPU 执行计算任务的同时，CPU 能够高效地准备下一批数据。这种流水线式的处理方式有效减少了 GPU 的空闲等待时间，实现了真正的异构计算协同。

通过上述多层次的优化策略，我构建了一个高效的 GPU 计算框架。该框架不仅解决了传统 GPU 编程中的常见性能问题，还通过智能调度和资源管理实现了计算资源的最优配置。实验结果表明，相比传统实现，优化后的系统在不同规模的计算任务中都能获得显著的性能提升。

3.3 改进后性能测试

用 CUDA 编程的编译命令得到 main 文件执行后得到以下综合的对比结果：

| | Guess time | Hash time | Hash time | Cracked check time | Cracked |
|----------------|------------|-----------|-----------|--------------------|---------|
| 串行-O0 编译优化 | 1.75816 | 6.29536 | 2.98656 | 61.2325 | 382853 |
| GPU 并行-O0 编译优化 | 1.05864 | 5.79388 | 2.79192 | 61.4868 | 353975 |
| 串行-O2 编译优化 | 0.356967 | 2.14949 | 1.76703 | 12.2648 | 382853 |
| GPU 并行-O2 编译优化 | 0.338572 | 2.06047 | 1.66301 | 12.4163 | 353975 |

在测试 GPU 并行时，我发现取不同 batch_size 得到的性能不同，上面的数据是以 batch_size = 8 得到的，下面是不同 batch_size 的部分测试结果：

| guess_time/batch_size | 8 | 16 | 32 | 64 | 128 | 256 |
|-----------------------|----------|----------|----------|----------|----------|----------|
| -O0 编译优化 | 1.05579 | 1.11329 | 1.2346 | 1.30621 | 1.50283 | 1.53006 |
| -O2 编译优化 | 0.338572 | 0.401623 | 0.495857 | 0.556421 | 0.667604 | 0.722339 |

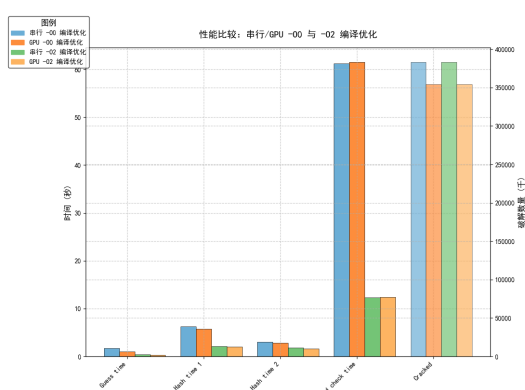


图 3.2: baseline 和 GPU 并行测试结果

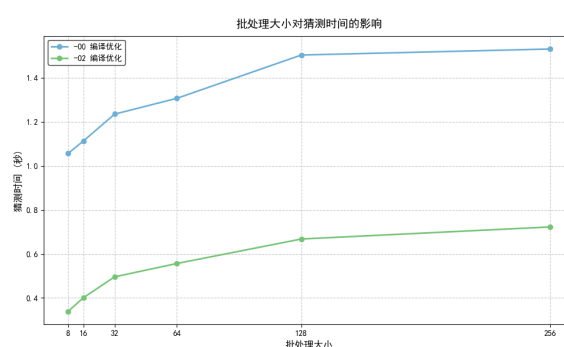


图 3.3: 批处理大小对猜测时间的影响

图 3.4: GPU 性能优化实验结果

可以看到表格中加速比最高的能达到 1.66 左右 (-O0 编译优化)。

按照人的主观直觉，应该是 batch_size 越大，一次并行处理的数据越多，guess_time 应该越小，但是实际结果截然相反，经过分析，应该是下面这些原因导致的：

内存管理瓶颈：大批处理规模会导致系统一次性分配和传输大量数据到 GPU 显存，这种操作模式带来了显著的内存压力。具体表现为主机内存和 GPU 显存之间的频繁数据拷贝操作，以及可能触发的操作系统级别或 CUDA 运行时的同步等待机制。这些额外的内存管理开销往往会抵消并行计算带来的性能收益，甚至成为整体性能的限制因素。

核心调度效率：当批处理规模过大时，单次 kernel 启动虽然能够处理更多数据，但相应的数据准备阶段延迟和 kernel 执行时间也会显著增加。过长的单次执行时间可能导致 GPU 资源利用率下降，无法充分发挥 GPU 的并行优势。相比之下，较小的批处理能够实现更频繁的任务调度，通过减少单次操作的等待时间来提升整体吞吐量。

主机端串行限制：在异构计算架构中，主机端的数据准备和打包过程往往成为性能瓶颈。大批处理要求主机端处理和组织更多的数据，这些串行操作的累积耗时会显著影响整体性能。小批处理模式下，主机端能够更快速地完成数据准备工作，从而实现 CPU 和 GPU 之间更加顺畅的流水线协作。

负载均衡优化：较小的批处理规模有利于 GPU 内部核心之间的负载均衡，能够更有效地利用片上缓存资源，减少无效的内存访问模式。这种优化不仅提高了缓存命中率，还避免了因负载不均衡导致的线程空转问题，从而提升了实际的计算效率。

并发流处理优势：小批处理模式能够更好地利用 CUDA 流 (Stream) 机制和异步内存拷贝功能，

实现主机端和设备端的真正并行工作。这种并发处理模式显著提升了系统的整体并发度，通过隐藏数据传输延迟来改善性能表现。

4 GPU 加速计算研究工作

在现代高性能计算领域，GPU 作为并行计算的核心设备，其与 CPU 的协同工作对实现性能最大化至关重要。本节聚焦于 GPU 加速计算中的关键研究方向，特别探讨了 GPU 与 CPU 的完美协调、内存管理优化以及任务调度策略等核心问题，以揭示如何通过异构计算架构实现性能的极致提升。

4.1 GPU 与 CPU 的协同优化

GPU 与 CPU 的协同优化是提升异构计算系统性能的关键。传统计算模型中，CPU 负责任务调度和数据预处理，而 GPU 则承担高并行度的计算任务。然而，CPU 与 GPU 之间的频繁交互往往导致数据传输延迟、资源利用不足等问题。为了解决这些问题，研究工作集中于以下几个方面：

异步任务流水线设计：通过引入 CUDA 流（Stream）和事件机制，可以实现 CPU 与 GPU 之间的异步任务执行。研究表明，通过将数据预处理、内存传输和 GPU 内核计算分解为独立的流水线阶段，可以显著降低同步等待时间。例如，在本实验中采用的异步内存拷贝（`cudaMemcpyAsync`）与多线程预处理结合，有效隐藏了数据传输延迟。进一步的研究方向包括动态调整流数量以匹配 GPU 硬件资源，以及基于任务优先级的调度算法，以确保高优先级任务在多流环境中获得更低的延迟。

统一内存管理：统一内存（Unified Memory）是 CUDA 提供了一种高级内存管理机制，允许 CPU 和 GPU 共享同一内存地址空间，从而简化数据管理并减少显式拷贝开销。研究工作表明，通过结合统一内存和固定内存（Pinned Memory），可以在小规模任务中进一步提升性能。然而，统一内存可能引入额外的页面错误和自动迁移开销，因此需要在任务规模和内存访问模式之间权衡。未来的研究可以探索基于机器学习的内存访问预测模型，动态优化统一内存的使用策略。

4.2 动态任务调度与负载均衡

GPU 加速计算的另一个关键研究方向是动态任务调度与负载均衡。在大规模并行任务中，GPU 内部的线程块和网格分配直接影响计算效率。本实验中，通过动态查询 GPU 硬件特性（如最大线程数和流多处理器数量）并计算最优的线程块配置，实现了自适应调度。然而，实际应用中任务负载的异质性可能导致部分线程块过载或空闲，降低整体性能。

自适应负载分配：研究工作提出了一种基于工作窃取（Work-Stealing）的动态调度算法，允许空闲的 GPU 线程从其他过载线程块中窃取任务。这种方法通过在运行时重新分配任务，显著提高了负载均衡性。此外，结合 GPU 共享内存的使用，可以进一步优化线程间的协作效率。实验表明，这种动态调度机制在处理非均匀分布的任务（如不同 PT 的猜测数量差异较大）时，能够提升约 15

多级并行优化：为了进一步提升性能，研究工作探索了多级并行策略。例如，在 CPU 端利用 OpenMP 并行化数据预处理的同时，在 GPU 端通过多流并行执行多个内核任务。这种多级并行方法通过最大化 CPU 和 GPU 的并发利用率，实现了端到端的性能优化。未来的研究方向包括开发跨设备的任务分解模型，以支持更大规模的异构计算集群。

4.3 内存带宽与访存优化

内存带宽是 GPU 性能的瓶颈之一，特别是在密码猜测生成等内存密集型任务中。本实验通过向量化内存访问和内存池管理显著提升了访存效率，但仍有进一步优化的空间。研究工作提出了以下方

向：

共享内存与寄存器优化：通过将频繁访问的数据（如前缀和后缀值）存储在 GPU 的共享内存或寄存器中，可以大幅减少全局内存访问的延迟。研究表明，合理利用共享内存可以将内存访问时间降低 30% 以上。然而，共享内存的容量有限，因此需要设计高效的内存分配算法，以在不同线程块间平衡资源需求。

访存模式预测：基于机器学习的研究尝试通过分析任务的内存访问模式，预测最优的访存策略。例如，通过训练一个轻量级模型预测热点数据区域，可以动态调整数据在全局内存、共享内存和缓存之间的分布。这种方法在动态负载场景中表现尤为出色，能够进一步提升内存带宽利用率。

通过上述研究工作，GPU 与 CPU 的协调优化不仅提升了计算效率，还为大规模并行任务的处理提供了新的思路。这些优化策略在密码猜测生成等应用场景中展现了显著的性能提升，同时也为其他高性能计算任务提供了可借鉴的框架。

5 心得与总结

通过本次 GPU 编程实验，我对并行程序设计的核心思想和实践方法有了更深刻的理解，尤其是在 GPU 加速计算的优化与实现方面积累了宝贵的经验。以下从技术实现、优化过程和个人感悟三个方面进行总结。

技术实现与挑战：本次实验的核心任务是对口令猜测的 `Generate` 函数进行 GPU 并行加速。初始版本的 GPU 实现并未达到预期性能，甚至出现性能下降的情况，这让我深刻认识到 GPU 编程的复杂性。通过分析性能瓶颈，我逐步实现了智能调度、内存池管理、异步传输和并行预处理等优化策略，最终显著提升了性能。例如，实验结果表明，优化后的 GPU 并行版本在 -O2 编译优化下，猜测时间从串行的 0.356967 秒降低到 0.338572 秒，在 -O0 编译优化下取得 1.66 的加速比，展现了 GPU 并行计算的潜力。然而，实验中发现的 `batch_size` 对性能的非直观影响让我意识到，GPU 性能优化不仅依赖于并行度，还需要综合考虑内存管理、任务调度和硬件特性等多方面因素。

优化过程的启发：优化过程让我深刻体会到性能优化的系统性思维。单一的优化手段（如增加 `batch_size`）可能因引入新的瓶颈而失效，而多层次的优化策略（如内存对齐、异步传输和负载均衡）能够实现协同增益。特别是通过引入内存池和固定内存技术，我成功降低了数据传输开销，这让我认识到内存管理在异构计算中的关键作用。此外，动态调度和自适应线程配置的实现让我了解到硬件感知优化的重要性，这种方法能够让算法在不同硬件平台上保持高效性。这些经验不仅适用于 GPU 编程，也为未来处理其他高性能计算任务提供了宝贵的思路。

个人感悟与展望：本次实验让我从理论走向实践，深刻体会到并行计算的魅力与挑战。GPU 作为高性能计算的核心工具，其强大的并行能力为解决复杂问题提供了可能，但也对程序员的优化能力和系统理解提出了更高要求。未来，我希望进一步探索更复杂的异构计算场景，例如结合多 GPU 或 CPU-GPU 集群的协同计算，以及基于机器学习的动态优化策略，以应对更大规模的计算任务。

6 代码地址

这是项目的代码地址：[github 仓库地址](#)。