



西北工业大学

本科毕业设计（论文）

题 目 基于华为鲲鹏 920 的 SpTRSV 算法优化

专业名称 计算机科学与技术

学生姓名 严愉程

指导教师 王云岚

完成时间 2021 年 6 月 1 日

毕业 设计 任务书

论文

一、题目

题目这种东西随便起一个就行了

二、研究主要内容

用 word 做好任务书, 打印成 pdf.

三、主要技术指标

50 页

四、进度和要求

第 1 周至第 2 周:完成第一章;
第 3 周至第 4 周:完成第二章;
第 5 周至第 6 周:完成第三章;
第 7 周至第 8 周:完成第四章;
第 9 周至第 10 周:完成第五章;
第 11 周至第 15 周:安心养老, 等待毕业;

五、主要参考书及参考资料

[1] Shen S. LaTeX-Template-For-NPU-Thesis[Z]. 2016.05.

学生学号 _____ 学生姓名 _____

指导教师 _____ 系 主 任 _____

摘 要

稀疏下三角矩阵求解器，是数值计算中的一个重要的方法，在科学计算中有着广泛的应用。由于其离散的数据访存、任务之间存在着很强的依赖、需要细粒度的同步、任务之间负载不均衡等特点，在现代多核系统中，如何提升 SpTRSV 算法的并行度与扩展性，仍然是一个挑战性的课题。

华为鲲鹏 920 处理器基于 ARMv8 架构，最多拥有 64 个核心，支持 ARM NEON 128bits 浮点运算。华为鲲鹏系统支持 NUMA 内存架构，实现最多 4 个 920 互联，最高支持 256 个计算核心。

本文基于华为鲲鹏 920 芯片，进行稀疏下三角矩阵求解算法的设计与优化。主要工作包括：

1. 设计并实现了一套高效的稀疏下三角矩阵求解算法。
2. 针对华为鲲鹏 920 处理器众核体系结构以及 ARMv8 指令集的特点进行了优化。
3. 采用 ARM NEON 指令，对向量乘法部分进行了 SIMD 的优化。
4. 结合处理器 NUMA 架构的特性，通过分配矩阵数据的存储，充分利用内存带宽。

关键词： 稀疏下三角矩阵求解器, 鲲鹏 920, 众核优化, 并行算法

Abstract

The sparse triangular solve kernel, SpTRSV, is an important building block for a number of numerical linear algebra routines. It is commonly considered challenging to parallelize and scale even on a moderate number of cores. This challenge is due to the fact that triangular solver typically has frequent discrete memory access, load imbalances between the tasks.

Huawei Kunpeng920 is an ARMv8-based server CPU with up to 64 cores, supporting ARM NEON 128-bit SIMD instruction set.

The main work of this article includes:

1. Design and implement an efficient sparse matrix triangular solve algorithm.
2. Optimizing for Huawei Kunpeng920 multi-core processor and features of ARMv8 instruction set.
3. Utilizing ARM NEON 128bits SIMD instruction set to speed up the vector multiplication which is part of SpTRSV algorithm.
4. By taking NUMA feature into consideration and allocating matrix data properly, achieve better use of memory bandwidth.

Key Words: SpTRSV, Kunpeng920, Multi-core Optimization, Parallel Algorithm

目 录

摘要	i
ABSTRACT(英文摘要)	ii
目录	iii
第一章 绪论	1
1.1 课题背景	1
1.2 研究现状	1
1.3 研究内容	3
1.4 本文构成	3
第二章 相关理论与技术介绍	4
2.1 线性方程组求解	4
2.2 稀疏矩阵的压缩方式	4
2.2.1 COO Format	4
2.2.2 CSR Format	4
2.2.3 CSC Format	5
2.3 现存 SpTRSV 算法	5
2.3.1 基于 CSR 格式的 SpTRSV 串行算法	5
2.3.2 基于 CSC 格式的 SpTRSV 串行算法	5
2.3.3 基于 level-sets 的 SpTRSV 并行算法	6
2.3.4 sync-free 的 SpTRSV 算法	7
2.4 基于缓存的优化技术	8
2.4.1 内存 cache	8
2.4.2 Cache Prefetch 优化技术	8
2.4.3 缓存一致性协议	10
2.4.4 伪共享	11
2.5 ARM Neon 指令	12
2.6 NUMA 架构	12
2.7 ARMv8 原子指令优化	13
2.8 本章小结	15
第三章 SpTRSV 算法的设计与实现	16
3.1 总结总结	16
3.2 总结总结	16
3.3 总结总结	16

西北工业大学 本科毕业设计论文

参考文献	17
附录	19
致谢	20
毕业设计小结	21

第一章 绪论

1.1 课题背景

稀疏下三角矩阵求解 (Sparse Triangular Solve, SpTRSV) 并行算法是很多数值计算方法的重要组成部分, 比如用直接法求解稀疏矩阵的线性方程组 [1], 以及最小二乘法的求解 [2], 是现代科学计算中一个广泛使用的计算核心, 在数值模拟计算中, 通常会使用迭代法或直接法求解大规模稀疏线性方程组, 而 SpTRSV 的效率直接影响了线性方程组的求解效率, 故提高 SpTRSV 算法的性能至关重要。

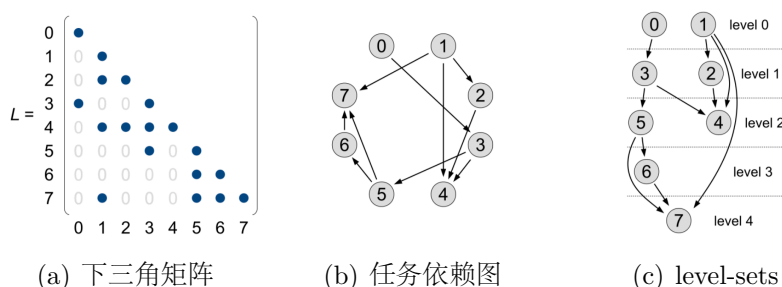
相比于稠密下三角举证解法器 [3] 或者其他稀疏矩阵计算方法, 例如稀疏矩阵转置 [4], 稀疏矩阵向量乘法 [5][6], 和稀疏矩阵乘法 [7], SpTRSV 频繁且离散地数据访存、任务之间存在着很强的依赖、需要细粒度的同步、任务之间负载不均衡等特点是并行优化更加难以进行。

SpTRSV 算法主要是从方程组 $Lx = b$ 中求解未知数向量 x , 其中 L 是下三角矩阵, 且对角线上的元素都不为 0。因为这意味着要计算其中的一个未知数 x_k , 需要先计算出其前置未知数 x_0, \dots, x_{k-1} 中的一个子集。这些任务间的依赖关系, 就使得 SpTRSV 的计算转换成了一个任务依赖图的计算 (Task Dependent Graph, TDG), 而且这个 TDG 是一张有方无环图 (Directed Acyclic Graph)。

华为作为鲲鹏计算产业的成员, 聚焦于发展华为鲲鹏 + 昇腾双引擎芯片族, 通过“硬件开放、软件开源、使能合作伙伴”来推动计算产业的发展。华为鲲鹏 920 处理器兼容了 ARMv8 架构。最多拥有 64 核心。其浮点及 SIMD 运算单元支持每拍 2 条 ARM Neon 128bits 浮点运算。华为鲲鹏系统支持 NUMA 内存架构, 实现最多 4 个鲲鹏 920 互联和最高 256 个物理核的 NUMA 架构。

稀疏矩阵运算在科学计算、数据分析、机器学习等应用中十分常见, 而稀疏矩阵频繁且离散地数据访存、任务之间存在着很强的依赖、需要细粒度的同步、任务之间负载不均衡等特点又给代码优化带来了巨大的挑战。

1.2 研究现状



有研究者 Saad[8], 以及 Saltz[9] 根据 SpTRSV 算法是一张任务依赖

图2-4(b)的特点,提出了基于 level-sets 的方法,在预处理阶段构建一个任务依赖图,如图1-1(c),之后就可以在 level 内部进行并行,在 level 之间设置同步屏障,一个 level 执行完了再执行下一个 level。

随后有研究人员在体系结构层面上针对 CPU[10][11][12][13] 以及 GPU[14][15][16] 进行了优化。

基于 level-sets 的算法,存在着两个缺陷:第一,虽然可以在每个 level 上进行并行取得良好的并行度,但是需要在预处理阶段构建一个 TAG 图,计算每个任务的 level,以及处理任务之间任务负载均衡的问题。然而,在预处理阶段往往会花费大量的时候,写出一个有良好扩展性的并行预处理算法同样具有挑战性,往往可能会导致预处理的时间远远大于并行计算任务图的时间,甚至在计算任务图上获得的加速甚至不能抵消预处理阶段产生的计算时间。第二,基于 level-sets 的的算法需要在每个 level 之间都要使用一个屏障确保进行同步,等待该 level 内的任务都执行完了再继续下一个任务,在负载不均衡的情况下,这会产生大量的等待时间,随着核心数量的上升该同步方式的开销也会随之上升。

面对以上两个 level-sets 的缺陷,有研究人员提出了以下的优化算法。

Jongsoo Park[11] 在基于 level-sets 算法进行了一些优化。他发现传统的基于 level-sets 的算法所使用的屏障式的同步方式会产生大量的开销。因此作者提出了一种 peer-to-peer 的同步方式。线程在任务执行完了之后不再是等待在屏障处,等待其他线程执行到该屏障处,而是会判断下一个任务的前置需求是否被完成,如果已经完成了,那么就继续运行下去。相比于基于屏障的同步方式,peer-to-peer 的同步方式具有更好的扩展性。

在减少同步所需消耗方面,作者还发现,在 SpTRSV 算法的任务依赖图中,大部分的依赖都是多余的,作者期望通过一步预处理操作来删除这些多余依赖。例如在 $2 \rightarrow 3 \rightarrow 5$ 和 $2 \rightarrow 5$ 这种边存在的情况下,删除了 $2 \rightarrow 5$ 这条边。另一方面,作者发现在这些多余依赖当中,大部分都是两跳的(比如上述的 $2 \rightarrow 5$ 这条边),少部分是三跳或者三跳以上。为了减少预处理的计算量,作者选择用粗略但快速的算法,只删除两跳的依赖边,而不是删除所有的多余依赖。减少了大约 90% 的多余依赖,具体效果如1-1。该算法运行在 12 核的 Xeon 处理器上,相比于传统的基于 level-sets 以及屏障式同步的算法获得了至少 1.6 倍的加速比。

同时作者也做了负载均衡的优化,根据每行非零元个数的多少来进行任务的合并,组合成一定数量的 super-task,使得每个 super-task 有相似的计算量。

刘伟峰 [17, 18] 从减少算法预处理和减少算法同步所需时间的角度出发,提出了无需同步的 (synchronization-free) SpTRSV 算法。基于传统 level-set 的方法,随着矩阵大小的增加,其在预处理阶段所需的消耗会剧烈增加,且同步所需的时间在计算总时间中的占比也会大量增加。他的算法在预处理阶段只计算了每个任务的 in-degree 数组(用来记录每一个节点需要有多少入度),而不是构造一个依赖图,大幅减少了预处理时间,而且该预处理方法能

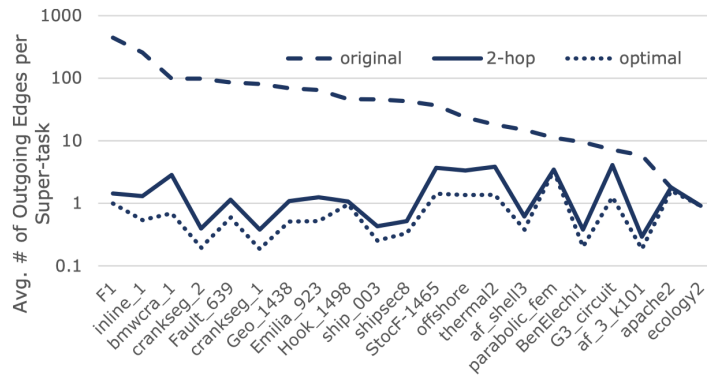


图 1-1 每个 super-task 的出度有了显著的减少

够轻易地使用多核来进行并行加速。对于每个 warp 使用自旋锁，不断判断前置依赖是否都被满足来决定该结点是否开始计算。当一行计算完成后，使用原子操作来“通知”后继节点，减少了同步的时间。同时作者结合 GPU 的体系结构，针对 GPU 的片上内存和全局内存进行了优化，该算法在 GPU TitanX 上比 Nvidia 提供的算法快 2-3 倍。

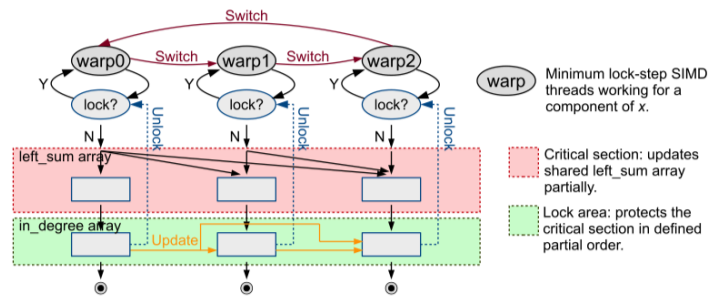


图 1-2 synchronization-free SpTRSV 算法的示意图

倪鸿、刘鑫 [19] 根据国产异构众核处理器 SW26010 体系结构的特点，针对非结构网络计算，提出了一种基于流水线串行-局部并行思想的通用众核 SpTRSV 优化方法。首先，为了减少预处理时间，根据核心的个数，将向量 x 平均划分为多个向量块。每个向量块 x 交给一个从核进行处理。每个从核的操作流程为：(1) 等待接受来自依赖向量块的数据，直到满足进入 (2)。(2) 进行 x 的计算。(3) 发送计算得到的 x 给需要的核心。其次，为了解决 SW26010 寄存器通讯的局限性，作者还搭建了众核通信架构。同时还有数据压缩、LDM 缓冲、访存掩盖以及数据压缩等优化。

1.3 研究内容

1.4 本文构成

第二章 相关理论与技术介绍

2.1 线性方程组求解

经典的求解线性方程组的方法一般分为两类：直接法和迭代法。前者例如高斯消元法, LU 分解等, 后者的例子包括共轭梯度法等。直接法指在不考虑计算舍入误差的情况下, 通过包括矩阵分解和三角方程组求解等有限步的操作求得方程组的精确解, 因此又称精确法; 迭代法指给定一个初始解向量, 通过一定的计算构造一个向量列 (一般通过逐次迭代得到一系列逼近精确值的近似解), 向量列的极限为方程组理论上的精确解。

2.2 稀疏矩阵的压缩方式

存储矩阵的一般方法是采用二维数组, 其优点是可以随机地访问每一个元素, 因而能够容易实现矩阵的各种运算。对于稀疏矩阵, 它通常具有很大的维度, 有时甚大到整个矩阵 (零元素) 占用了绝大部分内存采用二维数组的存储方法既浪费大量的存储单元来存放零元素, 又要在运算中浪费大量的时间来进行零元素的无效运算。因此必须考虑对稀疏矩阵进行压缩存储 (只存储非零元素)。

2.2.1 COO Format

最简单的稀疏矩阵的存储格式称为 coordinate(COO) format, 这个形式只保留哪些非 0 的值, 分别使用三个数组: val, rowIdx, colIdx 来对应存储数值、行号以及列号, 这三个数组的大小都为矩阵非零元的个数 NNZ。举例说明,

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 2 & 4 & 0 & 6 & 0 \\ 3 & 0 & 5 & 0 & 0 \\ 4 & 7 & 0 & 8 & 9 \end{pmatrix}$$

该矩阵在 COO 存储格式下, 三个数组分别为:

COO Format

```
1: val = {1, 2, 4, 6, 3, 5, 4, 7, 8, 9}
2: rowIdx = {0, 1, 1, 1, 2, 2, 3, 3, 3, 3}
3: colIdx = {0, 0, 1, 3, 0, 2, 0, 1, 3, 4}
```

2.2.2 CSR Format

CSR (Compressed Sparse Row) 是一种按行压缩的矩阵存储形式。分别使用三个数组: val, rowPtr, colIdx 来对应存储数值、行指针、以及列号。其中 val 和 colIdx 的数组大小为非零元的个数 NNZ, 而 rowPtr 数组的大小一般为 m+1, 其中 m 为矩阵行的个数。例如图2-1。

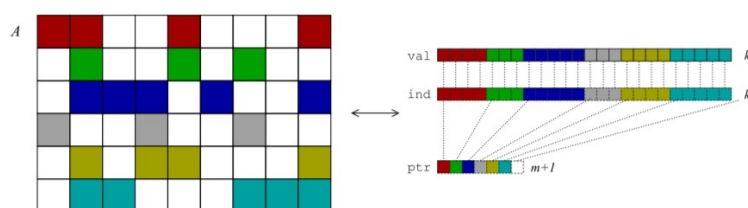


图 2-1 CSR 矩阵存储格式

2.2.3 CSC Format

CSC (Compressed Sparse Column) 与 CSR 相似，是一种按列压缩的矩阵存储格式。分别使用是那个数组：val, colPtr, rowIdx 来对应存储数值、列指针、以及行号。其中 val 和 rowIdx 的数组大小为非零元的个数 NNZ，而 colPtr 数组的大小一般为 $n+1$ ，其中 n 为矩阵列的个数。例如图2-2

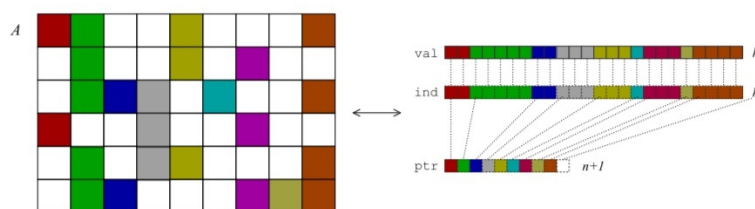


图 2-2 CSC 矩阵存储格式

2.3 现存 SpTRSV 算法

2.3.1 基于 CSR 格式的 SpTRSV 串行算法

Algorithm 1: CSR Based Serial SpTRSV

```

1  MALLOC(*left_sum, n) ;
2  MEMSET(*left_sum, 0) ;
3  for  $i = 0; i < n; i++$  do
4      for  $j = row\_ptr[i]; j < row\_ptr[i + 1] - 1; j++$  do
5           $left\_sum[i] \leftarrow left\_sum[i] + val[j] * x[col\_idx[j]]$  ;
6      end
7       $x[i] \leftarrow (b[i] - left\_sum[i]) / val[col\_ptr[i]]$  ;
8  end
    
```

基于 CSR 格式的 SpTRSV 串行算法第一个循环从矩阵的第 0 行开始遍历这个矩阵。在第二个循环进行向量乘法的操作，求出 $left_sum[i]$ 。之后再求出未知数 $x[i]$ 。

2.3.2 基于 CSC 格式的 SpTRSV 串行算法

该算法的第一个循环从矩阵的第 0 行开始遍历这个矩阵。首先求出该行未知数 $x[i]$ ，接着遍历该列的上的非零元，对该列非零元所对应的

Algorithm 2: CSC Based Serial SpTRSV

```

1  MALLOC(*left_sum, n) ;
2  MEMSET(*left_sum, 0) ;
3  for  $i = 0; i < n; i++$  do
4       $x[i] \leftarrow (b[i] - left\_sum[i]) / val[col\_ptr[i]]$  ;
5      for  $j = col\_ptr[i] + 1; j < col\_ptr[i + 1]; j++$  do
6           $left\_sum[rowIdx[j]] \leftarrow left\_sum[row\_idx[j]] + val[j] * x[i]$  ;
7      end
8  end

```

$left_sum$ 进行更新, $left_sum[rowIdx[j]] += val[j] \times x[i]$ 。

基于 CSR 压缩格式 SpTRSV 算法以及基于 CSC 压缩格式的 SpTRSV 算法, 两者在访存模式上存在着一定的差异。基于 CSR 格式的算法需要进行频繁离散读取未知数向量 $x[col_idx[j]]$, 相反基于 CSC 格式的 SpTRSV 算法只需要读取当前行号 i 所对应的 $x[i]$ 。在写数据方面, 基于 CSR 的并行算法需要集中的写 $left_sum[i]$ 的数据, 而基于 CSC 的并行算法则是离散的写 $left_sum[rowIdx[j]]$, 在多核运算的条件下, 离散的写比集中的写同一位置的数据更容易进行并行化。

2.3.3 基于 level-sets 的 SpTRSV 并行算法

J.Park 在 level-sets 算法的基础上提出除了更高效的并行 SpTRSV 算法 [11], 该算法主要分为两个阶段: 预处理阶段和计算阶段。

在预处理阶段主要有三个处理步骤:

1. 使用由 Chhugani[20] 提出的高度优化的并行广度优先搜索算法 (BFS), 构建了一个任务依赖图2-4(b)。
2. 在任务依赖图的基础上, 作者又通过将同 level 的任务合并为 super-task 的形式实现负载均衡, 确保每个线程被分配有相似计算量的 super-task, 这种合并的操作有进一步减少了任务之间的依赖个数, 进一步减少了同步所需的消耗。
3. 作者发现任务依赖图中存在着大量的多余依赖边, 并开发了一个高效的过滤算法, 去除了两条的依赖边。

在计算阶段作者使用了与传统 level-sets 算法不同的调度方式, 作者使用 peer-to-peer 的同步方式替换为了在 level 之间使用屏障3进行同步的方式。下面是两种不同同步方式的伪代码: 算法3、算法4

对比两种不同的调度算法, 可以发现使用 peer-to-peer 的调度方式, 能够使得计算资源能够得到更好的利用。因为当一个线程完成任务计算之后, 基于屏障的同步方式会让该线程进入等待的状态。而基于 peer-to-peer 的同步方式则会让该线程继续进行下去。如图2-3, 如果采用基于 peer-to-peer 的

Algorithm 3: Level-scheduling with barrier

```

1 foreach level l do
2   solve the unknowns of  $task_l^t$ ;
3   barrier;
4 end

```

Algorithm 4: Level-scheduling with peer-to-peer synchronization

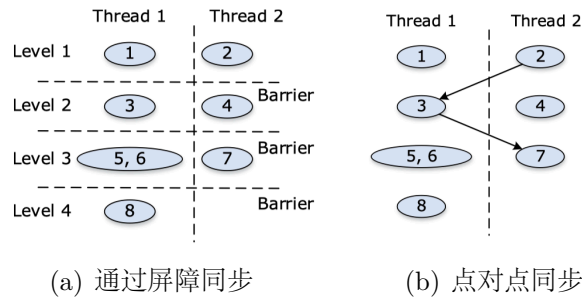
```

1 foreach level l do
2   wait until all the parents are done wait ;
3   solve the unknowns of  $task_l^t$ ;
4   done[ $task_l^t$ ];
5 end

```

同步方式，线程 2 在完成任务 2 之后可以不用等在线程 1 完成任务 1，可以接着运行下去。这在一定程度上也有利于负载均衡。

图 2-3 两种不同同步方式的对比示意图



2.3.4 sync-free 的 SpTRSV 算法

基于 level-sets 的算法需要花费大量的时间在预处理和任务间的同步上，运行 J.Park 的算法并分步进行统计（预处理阶段，同步所需耗时，计算阶段），我们就可以得到表2-1。所用的稀疏矩阵来自于佛罗里达大学的稀疏矩阵集合 [21]，观察这张表我们可以主要得到以下几点信息：

1. 预处理阶段所消耗的时间，要远大于计算阶段所消耗的时间，平均大概是 8 至 10 倍左右。预处理占总时间大约 80%。
2. 同步所需的开销会导致的时间消耗，会随着矩阵 level 个数的增加而剧烈增加，当稀疏矩阵的 level 较多时，会导致同步的时间消耗远大于计算的时间消耗。这也可能是同 level 间负载不均衡所导致的。

例如 nlpkkt160 只有两个 level，同步所占用的时间要远小于计算所需的时间。作为对比矩阵 FEM/ship 003，虽然矩阵的规模小于 lpkkt160，但是由于其 level 相对较多，导致了同步所占用的时间也较多。

Matrix name	Preprocessing cost (ms)	SpTRSV cost (ms)	SpTRSV cost breakdown (ms)		#Level-sets
			Synchronization	Compute	
FEM/ship 003	92.46	12.95	10.96	1.99	4367
FEM/Cantilever	47.89	9.60	5.62	3.98	2397
chipcool0	8.74	1.99	1.15	0.84	534
nlpkt160	484.67	38.30	0.01	38.29	2

表 2-1 基于 level-sets 算法 [11] 的任务耗时分解图

刘伟锋从减少预处理和减少同步所需时间的角度出发，提出了名为 sync-free 的 SpTRSV 计算方法 [17]。该算法能够在 GPU 上取得较好的效果。与传统基于 level-sets 的算法不同的是，sync-free 的 SpTRSV 算法只在预处理阶段计算每个任务的入度，该入度表示当前任务有多少的前置依赖。作者使用一个 warp，32 个线程来处理每个任务。任务之间通过使用共享变量以及原子操作的方式，来决定任务之间的顺序。算法的伪代码见 5

2.4 基于缓存的优化技术

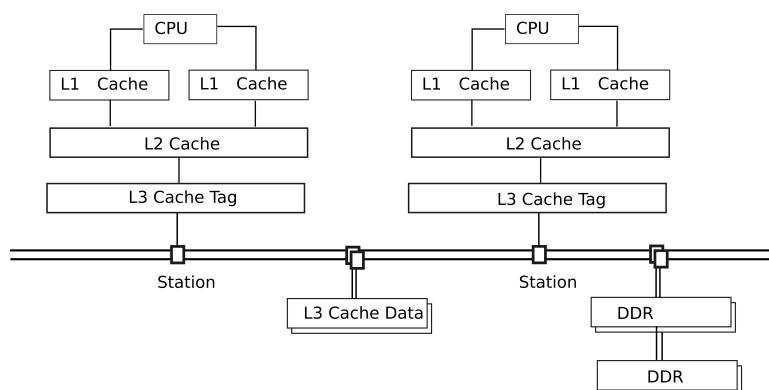


图 2-4 kunpeng cache 示意图

2.4.1 内存 cache

Cache 在计算机系统中，是一种基于分层优化思想的一种优化手段。由于 CPU 的速度远大于内存的速度，出于成本的考虑，人们常用相对告诉的 Cache 来弥补两者之间的差距，提升 CPU 的性能。在鲲鹏 920 芯片中，每个核心独占 L1、L2 两级缓存。离核心最近的是两个 L1 Cache，分开缓存指令和数据，大小分别为 64K。L2 Cache 则要大多得多，大小为 512k。而第三级缓存则是鲲鹏 920 芯片上的所有核心共享，其中 L3 Cacheline 的长度为 128 个字节。

2.4.2 Cache Prefetch 优化技术

Cache Prefetch 也是一个针对 Cache 的优化设计。Cache 比实际的内存快很多，所以如果我们可以提前加载部分内存到 Cache 中，就会在性能上有优势。比如一个乘法的时间，就算不考虑流水线，也不过 3-12 个周期，但是

Algorithm 5: sync-free SpTRSV

```

1  malloc(*d_left_sum, *s_left_sum, *d_in_degree, *s_in_degree, n);
2  memset(*d_left_sum, *s_left_sum, *d_in_degree, *s_in_degree, 0);
3  Function Preprocess() is
4      Parallel for  $i = 0; i < nnz; i++$  do
5          |   atomic-add(&d_in_degree[row_idx[i]], 1) ;
6      end
7  end
8  Function Solve() is
9      Parallel for  $int\ i = 0; i < n; i++$  do using one warp for one entry
10         |   while  $s\_in\_degree[i] + 1 \neq d\_in\_degree[i]$  do test dependencies
11         |   |   // busy wait ;
12         |   end
13         |    $x[i] \leftarrow (b[i] - d\_left\_sum[i] - s\_left\_sum[i]) / val[colptr[i]]$ ;
14         |   Parallel for  $j$  in  $[col\_ptr[i] + 1, col\_ptr[i + 1] - 1]$  do one thread for
           |   one nonzero
15         |   |   if in local then
16         |   |   |   atomic-add(&s_left_sum[rid],  $val[j] \times x[i]$ ) ;
17         |   |   |   atomic-add(&s_in_degree[rid], 1) ;
18         |   |   else
19         |   |   |   atomic-add(&s_left_sum[rid],  $val[j] \times x[i]$ ) ;
20         |   |   |   atomic-add(&s_in_degree[rid], 1) ;
21         |   |   end
22         |   end
23     end
24 end

```

从 DDR 内存中读取数据可能会花费 100 个时钟周期。所以最好还是在计算开始前，先让数据从内存加载到缓存当中，做到计算和缓存同步进行。

例如，如下代码：

```

1:   for (i=0; i<W; i++) {
2:       for(j=0; j<W; j++) {
3:           c[i][j] = 0;
4:           if (!(j%INT_PER_CACHELINE))
5:               __builtin_prefetch((const void *)&c[i][j+
                    INT_PER_CACHELINE], 1, 3);
6:           for (k=0; k<H; k++) {
7:               c[i][j] += a[i][k]*b[k][j];
8:               if (!j && !(k%INT_PER_CACHELINE))
9:                   __builtin_prefetch((const void *)&
                    a[i][k+INT_PER_CACHELINE], 0,
                    3);
10:              if (!i && !(j%INT_PER_CACHELINE))
11:                  __builtin_prefetch((const void *)&
                    b[k][j+INT_PER_CACHELINE], 0,
                    3);
12:          }
13:      }
14:  }
    
```

这里的 `__builtin_prefetch` 是 gcc 的内置函数，在不同的平台有不同的封装。在上面的算法中，我们进行某个向量单元的计算的时候，已经可以知道后面要计算的内存单元是什么了，我们就可以提前把数据取出来。除了软件在做 prefetch 外，鲲鹏 920 处理器中也有硬件实现的 prefetch 单元，如 2-5。这就导致软件层面的 prefetch 可以没有很好的效果。

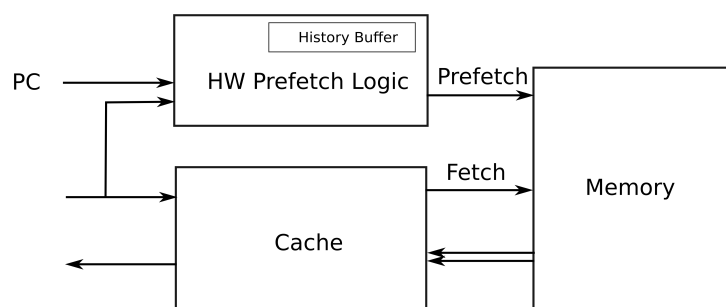


图 2-5 鲲鹏 920 cache prefetch 的硬件结构示意图

2.4.3 缓存一致性协议

在多核 CPU 的情况下，由于每个核心都有自己的独占的缓存，数据在这些缓存中制造了多个备份，这就造成了，到底哪份有效的问题。鲲鹏 920 采用的是监听一致性协议（Snooping），该协议是通过总线广播的形式实现

的。最为经典的监听一致性协议是 MESI，该协议由 James Goodman 提出，目前已在 x86、ARM、Power 架构中得到实现。

MESI 协议将 CacheLine 分为四个状态

1. Modified (M): 该状态表示，此 CacheLine 有效，数据只存在本 cache 中，且 CacheLine 中的数据被修改和内存中的不一样。在这种状态下，如果监听到有其他芯片读取缓存行对应主存的操作，该操作就会被延迟：首先将该操作写回，然后本 CacheLine 的状态就会变为 S。
2. exclusive (E): 该状态表示，此 CacheLine 有效，数据只存在于本 cache 中，数据和 Cache 中的数据和对应内存中的数据一致。如果监听到总线中其他线程对该行的读取操作，状态将变为 S。
3. Shared (S): 该状态表示，此 CacheLine 有效，数据存在于多个 Cache 中。如果监听到有其他线程的修改或独享请求，那么改缓存行就变为无效
4. Invalid (I): 该状态表示，此 CacheLine 无效，如果有线程对该数据进行读取，就触发 cache miss，重新从内存中读取对应的数据。

此外还有一种基于目录是的 MESIF 缓存一致性协议，这里就不详细展开，详见 [22]。

总线本质上是一个去中心化的系统，随着核心数量的增加，维护缓存一致性的成本也会增加。例如，假设 32 个线程同时在一个变量上做自旋锁的操作，此时如果有一个核心更新了一次 spinlock，那么就要通知另外 31 个线程。如果总线上发生了冲突，会导致计算性能进一步下降。

2.4.4 伪共享

缓存一致性协议以及数据按照 CacheLine 大小读取到 Cache，这两个特点造成了在多核计算中经常出现的问题：伪共享（false sharing）。

```
1:  #omp parallel for num_threads(10)
2:  for(int i = 0; i < 10; i++){
3:      for(int j = 0; j < 10; i++){
4:          sum[i] += i * j;
5:      }
6:  }
```

例如上述代码，在鲲鹏 920 中，CacheLine 的大小为 128 字节，sum 数组可能会被分配掉同一个 CacheLine 中，这种情况下，会导致算法的并行度无法上升。因为一个核心的写操作都会导致其他核心的缓存失效。当其他线程需要读写该区域内存的时候需要首先将修改过的 CacheLine 写回，然后重新从内存中读取。

通常可以通过 padding 的方式来解决伪共享的问题，如图2-6。对于多线程频繁读写的数据，通过 padding，将不想管的数据分配到不同的 CacheLine 当中。

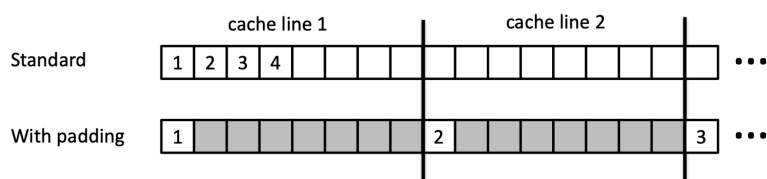


图 2-6 通过 padding 的方式解决伪共享导致的问题

2.5 ARM Neon 指令

鲲鹏 920 处理器基于 ARMv8，支持 128bits 的 SIMD 指令。它通过在 64bits 的 D 寄存器和 128bits 的 Q 寄存器上进行向量化的操作，来实现单指令多数数据运算。鲲鹏 920 处理器一共有 16 个 Q 寄存器和 32 个 D 寄存器。图2-7显示了指令 *VADD.I16Q0, Q1, Q2* 同时进行了 8 路 16bit 整数的加法操作，相比于传统的加法指令有显著的性能提升。ARM Neon 除了提供汇编指

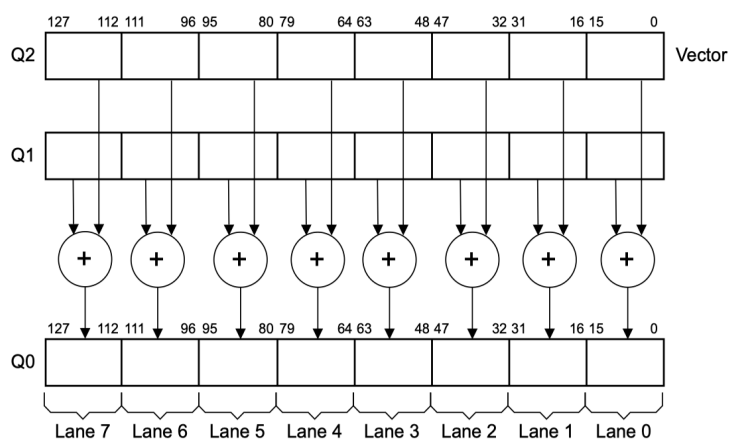


图 2-7 8 路 16 位加法操作

令之外，还提供了 C 风格的编程接口 Neon intrinsics。相比与 Neon 汇编指令，使用 Neon intrinsics 可以不用考虑寄存器的分配问题，这部分由编译器来实现，用户只需要关注上层代码。同时，也能享受编译器带来的一些优化。Neon intrinsics 了丰富的编程接口，包括：对数据进行的 load 与 store，各种常用的运算，变量常量的创建等。Neon intrinsics 支持 GCC 编译器，使用时只需 include 头文件 `arm_neon.h`。

2.6 NUMA 架构

传统多核处理器通常采用 SMP (Symmetric Multi-Processing, 对称多处理器) 2-8，在该架构下，每个线程的地位都是均等的，对内存使用的延迟与带宽也相同。在操作系统的支持下能够做到很好的负载均衡。鲲鹏处理器支持 NUMA (Non-uniform memory access, 非统一内存访问) 架构，如图2-9。使用这种架构能够解决处理器核数限制的问题，同时，通过合理的软件设计，能够提升内存的带宽，解决 SMP 架构下总线瓶颈的问题。

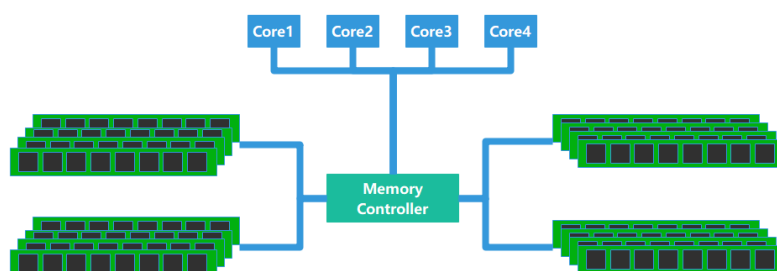


图 2-8 SMP 架构示意图

由于在 NUMA 架构下，多个核心组成一个节点 (Node)，系统中可以存在着多个节点，每个节点上有一个内存控制器，管理控制一片内存。内存在物理上是分布式的，通过片间网络互联。这就导致了，每个线程访问内存的性能取决于软件相对于内存的位置，对软件设计者提出了挑战。

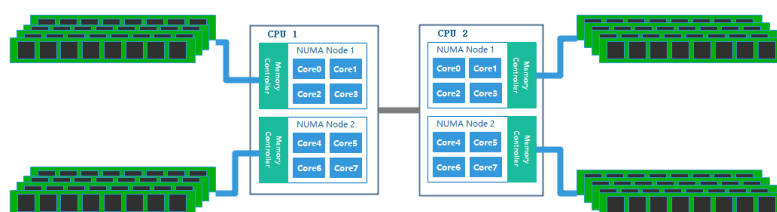


图 2-9 NUMA 架构示意图

linux 从内核版本 2.5 开始就支持 NUMA 架构 [23]，用户可以通过 numactl 命令来手动分配程序的 NUMA 访问策略，例如我们可以将进程或者内存固定在指定的节点上。系统默认的 NUMA 策略是 nodelocal，及总是在线程执行的地方分配内存。用户可以设置为 numactl -interleave，这时操作系统将会在各个节点之间，使用 round-robin 算法均衡分配内存。同时系统也为用户提供了 libnuma 这一编程接口 [24]，方便用户在程序语言层面，对数据进行操作，例如手动指定数据所在的节点。

值得一提的是，默认 NUMA 策略下，操作系统使用的是“Fist touch”的 NUMA 分配原则，即用户使用 malloc 时操作系统会分配逻辑内存页，当初次访问这个数据时（例如初始化操作），操作系统发现该逻辑页面还没有被分配物理页面，于是就进行分配物理页面的操作，此时操作系统的根据调用线程的 NUMA 分配策略来进行物理页面的分配。这时，在默认 NUMA 分配政策下，操作系统会将数据分配到，离该线程最接近的 NUMA 节点。

常用的编程工具 OpenMP 也提供了 NUMA 编程的接口。主要是通过两个环境变量，OMP_PLACES 决定了线程能够创建的位置，可以指定具体的节点，也可指定具体的核心。OMP_PROC_BIND 可以设置线程的分配方式。

2.7 ARMv8 原子指令优化

原子操作是多核程序中的一个很重要的操作。x86 通常使用锁总线的方式实现原子指令，保证了原子指令的执行期间不会收到其他指令的干扰。

在 ARMv8.1 发布之前，ARM 处理器实现原子操作的方式是 ll/sc(Load-Link/Store-Conditional)，具体指令就是 LDREX/STREX。

具体原理为，当 CPU-A 调用 LOADEX 命令的时候，会将数据从内存中读取到缓存当中，并将该缓存的状态标记为 EX (Exclusive)，此时如果还有一个 CPU-B 也使用了 LOADEX 命令，读取同一个数据，这时 CPU-B 中的标记变成了 EX，而 CPU-A 中缓存行的标记则变成了 S。在 CPU-A 执行完了之后如果要调用 STREX 指令，将数据写回内存，它首先检查标记的情况，发现自己没了 EX，于是写回失败。重新执行 LOADEX，操作数据，STREX 这样的循环。而 CPU-B 则可以成功进行 STREX，然后放弃 EX 标记。CPU-A 不断处于循环状态，这样就会导致性能的损失。

随着 ARM 处理器核心数量的提升，原先 LL/SC 方式的原子操作通过争抢的方式获得原子操作会对性能造成很大的影响。因此，为了支持这种大型系统，在其 ARMv8.1 规范中引入了 LSE (Large System Extensions)，加入了大量原生原子操作指令 [25]，大致原理就是将原子操作放到存储端去做，从而提升多核计算性能，理论上在多核系统的条件下，LSE 指令的性能要好与 ll/sc 原子指令。新的扩展指令包含 CAS, SWP 和 LD/ST<OP> 等，其中 <op> 为 ADD, CLR, SET 等。

gcc 也充分利用了这一特性，提供了基于 LSE 原子指令的优化。

CODE 2.1 基于 LL/SC 的原子指令

```
1:  __LL_SC_PREFIX(arch_atomic_##op(int i, atomic_t *v))
2:  {
3:      unsigned long tmp;
4:      int result;
5:      asm volatile("// atomic_" #op "\n"
6:      " prfm      pstl1strm, %2\n"
7:      "1: ldxr     %w0, %2\n"
8:      " " #asm_op " %w0, %w0, %w3\n"
9:      " stxr      %w1, %w0, %2\n"
10:     " cbnz      %w1, 1b"
11:     "=&r" (result), "=&r" (tmp), "+Q" (v->counter)
12:     : "Ir" (i));
13: }
```

CODE 2.2 基于 LSE 的原子指令

```
1:  static inline void arch_atomic_##op(int i, atomic_t *v)
2:  {
3:      register int w0 asm ("w0") = i;
4:      register atomic_t *x1 asm ("x1") = v;
5:
6:      asm volatile(ARM64_LSE_ATOMIC_INSN(__LL_SC_ATOMIC(op
7:      ),
```

```
7:      "      #asm_op " %w[i], %[v]\n")
8:      : [i] "+r" (w0), [v] "+Q" (v->counter)
9:      : "r" (x1)
10:     : "x16", "x17", "x30");
11:    }
```

对比上述两个不同的实现方式，原来基于 LL/SC 实现的原子操作，需要先 LDXR，然后执行 STXR 并查看是否成功，如果不成功那就重新循环执行，直到成功。而基于 LSE 原子指令的实现则更加简洁，如果是 atomic_add 操作只需一条 LDADD 指令。

2.8 本章小结

本章介绍了 SpTRSV 算法的应用场景，即线性方程组的求解；然后介绍了稀疏矩阵常见的两种几种压缩方式 COO、CSC、CSR；之后又详细介绍了目前性能不错的 SpTRSV 算法及其性能分析。最后则介绍了，基于多核 CPU 体系结构进行优化的一些手段，以及这些优化技术的原理。

第 三 章 SpTRSV 算法的设计与实现

3.1 总结总结

3.2 总结总结

3.3 总结总结

参考文献

- [1] Davis T A. Direct methods for sparse linear systems[M]. SIAM, 2006.
- [2] Saad Y. Iterative methods for sparse linear systems[M]. SIAM, 2003.
- [3] Hogg J D. A fast dense triangular solve in cuda[J]. SIAM Journal on Scientific Computing, 2013, 35(3):C303-C322.
- [4] Wang H, Liu W, Hou K, et al. Parallel transposition of sparse data structures[C]//Proceedings of the 2016 International Conference on Supercomputing. 2016: 1-13.
- [5] Liu W, Vinter B. Csr5: An efficient storage format for cross-platform sparse matrix-vector multiplication[C]//Proceedings of the 29th ACM on International Conference on Supercomputing. 2015: 339-350.
- [6] Liu W, Vinter B. Speculative segmented sum for sparse matrix-vector multiplication on heterogeneous processors[J]. Parallel Computing, 2015, 49:179-193.
- [7] Liu W, Vinter B. A framework for general sparse matrix-matrix multiplication on gpus and heterogeneous processors[J]. Journal of Parallel and Distributed Computing, 2015, 85:47-61.
- [8] Anderson E, Saad Y. Solving sparse triangular linear systems on parallel computers[J]. International Journal of High Speed Computing, 1989, 1(01):73-95.
- [9] Saltz J H. Aggregation methods for solving sparse triangular systems on multiprocessors[J]. SIAM journal on scientific and statistical computing, 1990, 11(1):123-144.
- [10] Kabir H, Booth J D, Aupy G, et al. Sts-k: a multilevel sparse triangular solution scheme for numa multicores[C]//Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. 2015: 1-11.
- [11] Park J, Smelyanskiy M, Sundaram N, et al. Sparsifying synchronization for high-performance shared-memory sparse triangular solver[C]//International Supercomputing Conference. Springer, 2014: 124-140.
- [12] Schreiber T, R. Vectorizing the conjugate gradient method.[J]. Proceedings of the Symposium on CYBER 205 Applications, 1982.
- [13] Wolf M M, Heroux M A, Boman E G. Factors impacting performance of multithreaded sparse triangular solve[C]//International Conference on High Performance Computing for Computational Science. Springer, 2010: 32-44.
- [14] Li R, Saad Y. Gpu-accelerated preconditioned iterative linear solvers[J]. The Journal of Supercomputing, 2013, 63(2):443-466.
- [15] Naumov M. Parallel solution of sparse triangular linear systems in the preconditioned iterative methods on the gpu[J]. NVIDIA Corp., Westford, MA, USA, Tech. Rep. NVR-2011, 2011, 1.
- [16] Suchoski B, Severn C, Shantharam M, et al. Adapting sparse triangular solution to gpus[C]//2012 41st International Conference on Parallel Processing Workshops. IEEE, 2012: 140-148.
- [17] Liu W, Li A, Hogg J, et al. A synchronization free algorithm for parallel sparse triangular solves[M]//Euro Par2016: Parallel Processing: volume 9833. ChamSpringer International Publishing, 2016: 617-630.
- [18] Liu W, Li A, Hogg J D, et al. Fast synchronization-free algorithms for parallel sparse triangular solves with multiple right-hand sides: Fast synchronization-free algorithms for parallel sparse triangular solves[J/OL]. Concurrency and Computation: Practice and Experience, 2017, 29(21):e4244. DOI: 10.1002/cpe.4244.

- [19] 倪鸿刘鑫. 非结构网格下稀疏下三角方程求解器众核优化技术研究[M]. 北京大学, 2019.
- [20] Chhugani J, Satish N, Kim C, et al. Fast and Efficient Graph Traversal Algorithm for CPUs: Maximizing Single-Node Efficiency[C/OL]//2012 IEEE 26th International Parallel and Distributed Processing Symposium. Shanghai, ChinaIEEE, 2012: 378-389. DOI: 10.1109/IPDPS.2012.43.
- [21] Davis T A, Hu Y. The university of florida sparse matrix collection[J]. ACM Transactions on Mathematical Software (TOMS), 2011, 38(1):1-25.
- [22] 胡森森计卫星. 片上多核处理器 Cache 一致性协议优化研究综述[M]. 软件学报, 2017.
- [23] Sarcar K. What is numa?[EB/OL]. <https://www.kernel.org/doc/html/v4.18/vm/numa.html>.
- [24] Kerrisk M. numa(3) —linux manual page[EB/OL]. <https://man7.org/linux/man-pages/man3/numa.3.html>.
- [25] Arm Architecture Reference Manual Armv8, for Armv8-A architecture profile[J]. :8248.

附 录

这是一份附录，请放置一些独立的证明、源代码、或其他辅助资料。

致 谢

感谢党和国家

感谢老师对我的指导和帮助

感谢中科院软件研究所和华为提供的计算资源

毕业设计小结

毕业论文是大学四年的最后一份大作业...