



南開大學  
Nankai University

计算机学院  
并行程序设计期末报告

并行加速的 SIMD 算法

姓名：刘迪乘  
学号：2312810  
专业：计算机科学与技术

2025 年 4 月 27 日

# 目录

<b>1 实验内容</b>	<b>2</b>
<b>2 实验环境</b>	<b>2</b>
<b>3 基础要求</b>	<b>2</b>
3.1 算法设计 . . . . .	3
3.2 代码实现 . . . . .	3
3.3 性能测试 . . . . .	5
3.4 结果分析 . . . . .	7
<b>4 进阶要求</b>	<b>7</b>
4.1 在 x86 架构下利用 SSE 指令集实现 SIMD 并行算法 . . . . .	7
4.2 尝试不同并行度 . . . . .	7
4.3 Profiling + 加速成功分析 . . . . .	8

## 1 实验内容

大问题：利用并行化手段，加速 PCFG 算法口令猜测的训练，猜测过程，同时对生成的猜测进行 MD5 算法加密。

本次问题：利用 SIMD 对 PCFG 算法生成的口令 md5 加密过程进行加速。

## 2 实验环境

本次实验选取的操作系统为 Windows11 下的 WSL(Linux 的版本为 24.04)，以及 windows 系统 (用于 X86)。此外作业提交在 OpenEuler 服务器 (ARM 架构) 上

硬件参数如下：

硬件名称 \ 参数	naive-conv
Architecture	x86_64
CPU(s)	16
Thread(s) per core	2
Core(s) per socket	8
Socket(s)	1
CPU max MHz	5000
CPU min MHz	3800
L1 cache	64K
L2 cache	8MB
L3 cache	32MB

表 1: AMD R7 9700X 硬件参数

硬件名称 \ 参数	naive-conv
Architecture	aarch64
CPU(s)	96
Thread(s) per core	1
Core(s) per socket	48
Socket(s)	2
CPU max MHz	2600
CPU min MHz	200
L1 cache	128K
L2 cache	512K
L3 cache	49152K

表 2: OpenEuler 服务器硬件参数

编译器版本：

- gcc version 13.3.0 (Ubuntu 13.3.0-6ubuntu2 24.04)
- gcc version 14.2.0 (MinGW-W64 x86\_64-ucrt-posix-seh, built by Brecht Sanders, r3)
- gcc (GCC) 9.3.1 20200408 (Red Hat 9.3.1-2)

代码开源地址:<https://github.com/July-h5kf3/NKU-parallel-computing/tree/main>

## 3 基础要求

尝试在 ARM 服务器上，基于 NEON 实现 MD5 哈希算法。

具体要求如下：

- 要求有基本的正确性，即能够利用 SIMD 指令，做到一次性产生多个消息（口令）的哈希值，并且哈希值正确。
- 不要求真正实现相对串行算法的加速，但是仍需要给出实验数据，并且分析为什么没有能够实现加速。

### 3.1 算法设计

并行的思路是将输入的消息预处理后涉及的 FF, GG, HH, II, 四个函数进行并行。这些函数并不涉及条件判断, 并且为较简单的运算 (移位, 与运算等)。因此可以让这些函数同时处理多个口令。从而实现并行化。

简单来讲, 首先将传入的多个口令预处理, 然后在分块处理时将每个口令 16 个 32bit 子块向量化 (uint32x4\_t), 再将各个函数对应运算操作替换为 ARM NEON 中对应的向量操作指令。具体代码见代码实现部分。

具体流程图如下, 并行加速的模块也有所体现:

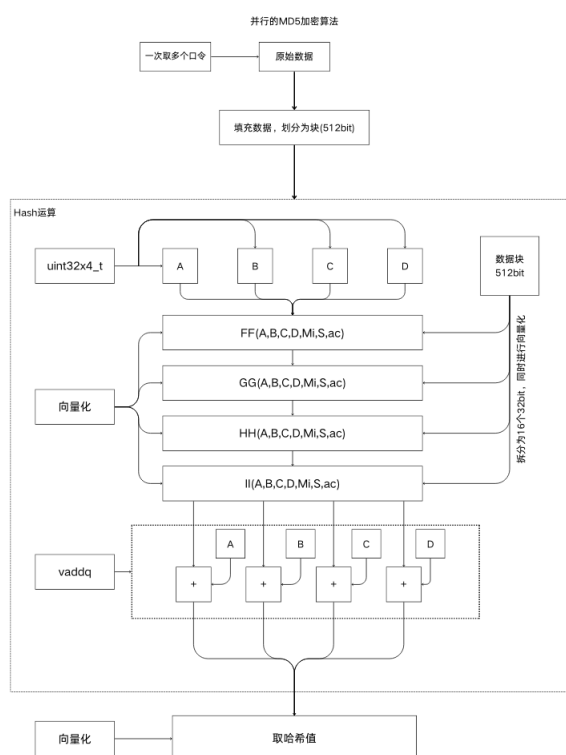


图 3.1: 流程图

### 3.2 代码实现

我从 main 函数接口部分着手, 首先将传入的参数由单个口令改为 4 个口令, 并适配相应的输出, 也就是把 256bit 的缓冲区进行向量化。

```

1 uint32x4_t state[4];
2 /*
3 这部分在hash函数中执行
4 state[0] = vdupq_n_u32(0x67452301);
5 state[1] = vdupq_n_u32(0xefcdab89);
6 state[2] = vdupq_n_u32(0x98badcfe);
7 state[3] = vdupq_n_u32(0x10325476);
8 */

```

```

9  int num = q.guesses.size();
10 for (int i = 0; i < num; i += 4)
11 {
12     string pw1 = q.guesses[i];
13     string pw2 = i+1 >= num ? "" : q.guesses[i + 1];
14     string pw3 = i+2 >= num ? "" : q.guesses[i + 2];
15     string pw4 = i+3 >= num ? "" : q.guesses[i + 3];
16     string pw_batch[4] = {pw1,pw2,pw3,pw4};
17     MD5Hash_SIMD(pw_batch, state);
18 }

```

然后，将各个口令进行初步的预处理，再将预处理后的口令逐 block 处理，由于口令长度的限制，各个口令的 paddedMessage 长度是相等的。对于每个 block，我首先将这个大小为 512bit 的 block 划分为 16 个 32bit 的子块。由于一次处理多个口令，因此要将多个口令得到的子块向量化：

```

1  uint32x4_t x[16];
2  bit32 y[4];
3  for(int i1 = 0; i1 < 16; i1++)
4  {
5      for(int j = 0; j < 4; j++)
6      {
7          y[j] = (paddedMessage[j][4 * i1 + i * 64]) |
8                (paddedMessage[j][4 * i1 + 1 + i * 64] << 8) |
9                (paddedMessage[j][4 * i1 + 2 + i * 64] << 16) |
10               (paddedMessage[j][4 * i1 + 3 + i * 64] << 24);
11      }
12      x[i1] = vld1q_u32(&y[0]);
13  }

```

之后就开始执行一系列的 FF，GG，HH，II 操作，下面展示并行化的相关函数：

#### md5 加速的核心部分

```

1  inline uint32x4_t ROTATELEFT_simd(uint32x4_t num, int n)
2  {
3      return vorrq_u32(vshlq_n_u32(num, n), vshrq_n_u32(num, 32-n));
4  }
5  inline uint32x4_t F_simd(uint32x4_t x, uint32x4_t y, uint32x4_t z)
6  {
7      return vorrq_u32(vandq_u32(x, y), vandq_u32(vmvnq_u32(x), z));
8  }
9  inline uint32x4_t G_simd(uint32x4_t x, uint32x4_t y, uint32x4_t z)
10 {
11     return vorrq_u32(vandq_u32(x, z), vandq_u32(y, vmvnq_u32(z)));
12 }
13 inline uint32x4_t H_simd(uint32x4_t x, uint32x4_t y, uint32x4_t z)
14 {
15     return veorq_u32(veorq_u32(x, y), z);
16 }
17 inline uint32x4_t I_simd(uint32x4_t x, uint32x4_t y, uint32x4_t z)

```

```

18 {
19     return veorq_u32(y, vorrq_u32(x, vmvnq_u32(z)));
20 }
21 // 考虑到FF,GG,HH,II 仅在调用F,G,H,I上存在差异, 因此只展示FF
22 inline void FF_simd(uint32x4_t *a, uint32x4_t b, uint32x4_t c, uint32x4_t d, uint32x4_t
    x, int s, uint32_t ac)
23 {
24     uint32x4_t ac_vec = vdupq_n_u32(ac);
25     *a = vaddq_u32(*a, vaddq_u32(F_simd(b, c, d), vaddq_u32(x, ac_vec)));
26     *a = ROTATELEFT_simd(*a, s);
27     *a = vaddq_u32(*a, b);
28 }

```

经过 4 轮操作后, 将缓冲区更新

```

1 state[0] = vaddq_u32(state[0], a);
2 state[1] = vaddq_u32(state[1], b);
3 state[2] = vaddq_u32(state[2], c);
4 state[3] = vaddq_u32(state[3], d);

```

最后对缓冲区进行移位操作

```

1 for(int i = 0; i < 4; i++)
2 {
3     uint32x4_t value = state[i];
4     uint32x4_t tmp1 = vdupq_n_u32(0xff);
5     uint32x4_t tmp2 = vdupq_n_u32(0xff00);
6     uint32x4_t tmp3 = vdupq_n_u32(0xff0000);
7     uint32x4_t tmp4 = vdupq_n_u32(0xff000000);
8     state[i] = vorrq_u32(vorrq_u32(vshlq_n_u32(vandq_u32(value, tmp1), 24),
9     vshlq_n_u32(vandq_u32(value, tmp2), 8)),
10    vorrq_u32(vshrq_n_u32(vandq_u32(value, tmp3), 8)
11    , vshrq_n_u32(vandq_u32(value, tmp4), 24)));
12 }

```

### 3.3 性能测试

首先展示我程序的正确性。

为验证程序的正确性, 我自行编写了如下代码

```

1 int main()
2 {
3     string pwd[4] = {"abc", "def", "hij", "knb"};
4     uint32x4_t state_simd[4];
5     //auto start_hash = system_clock::now();
6     MD5Hash_SIMD(pwd, state_simd);
7     //auto end_hash = system_clock::now();
8     //cout<<end_hash - start_hash<<endl;
9     bit32 state_show[4][4];

```

```

10     for (int i = 0; i < 4; i++) {
11         uint32_t tmp0[4], tmp1[4], tmp2[4], tmp3[4];
12         vst1q_u32(tmp0, state_simd[0]);
13         vst1q_u32(tmp1, state_simd[1]);
14         vst1q_u32(tmp2, state_simd[2]);
15         vst1q_u32(tmp3, state_simd[3]);
16
17         state_show[i][0] = tmp0[i];
18         state_show[i][1] = tmp1[i];
19         state_show[i][2] = tmp2[i];
20         state_show[i][3] = tmp3[i];
21     }
22
23     for(int i = 0; i < 4; i++)
24     {
25         bit32 state[4];
26         MD5Hash(pwd[i], state);
27         cout<<"correct:";
28         for(int j = 0; j < 4; j++)
29         {
30             cout << std::setw(8) << std::setfill('0') << hex << state[j];
31         }
32         puts("");
33         cout<<"me:";
34         for(int j = 0; j < 4; j++)
35         {
36             cout << std::setw(8) << std::setfill('0') << hex << state_show[i][j];
37         }
38         puts("");
39     }
40     return 0;
41 }

```

经验证，结果正确

对于我 SIMD 的性能评估，我主要从时间开销以及 profiling 等角度展开。

对于时间开销，我分别统计在不同优化力度下的串行、并行完成 PCFG 产生的所有口令 MD5Hash 过程的时间开销，以及完成不同数据规模的时间平均开销。

	SIMD	串行
无编译优化	14.2729s	11.477s
O1	2.57492s	5.49482s
O2	2.35317s	5.49505s
O3	2.14581s	5.31903s

表 3: 不同优化力度下的时间开销

加速比为**2.48**。对于不同数据规模的情况，由于服务器限制，选择在 X86 平台下对其他指令集的 SIMD 进行测试，测试结果见下面部分。

### 3.4 结果分析

## 4 进阶要求

### 4.1 在 x86 架构下利用 SSE 指令集实现 SIMD 并行算法

SSE 指令集下实现 SIMD 算法大体与 NEON 相同，只需要更改一些指令的名称，以及数据类型。具体而言，其 128bit 寄存器名称为 `__int128`。涉及的运算变更如下：

```

1  x[i1] = vld1q_u32(&y[0]); -> x[i1] =
    __mm_loadu_si128((__m128i*)y); // 从内存中加载数
2  state[0] = vdupq_n_u32(0x67452301); -> state[0] =
    __mm_set1_epi32(0x67452301); // 向量广播操作
3  vshlq_n_u32(num,n); -> __mm_slli_epi32(num,n); // 左移操作
4  vshrq_n_u32(num,n); -> __mm_srli_epi32(num,n); // 右移操作
5  vorrq_u32(x,y); -> __mm_or_si128(x,y); // 按位或操作
6  vandq_u32(x,y); -> __mm_and_si128(x,y); // 按位与操作
7  veorq_u32(x,y); -> __mm_xor_si128(x,y); // 按位异或操作
8  vaddq_u32(x,y); -> __mm_add_epi32(x,y); // 加法
9  vmvnq_u32(x) -> __mm_xor_si128(x,__mm_set1_epi32(0xFFFFFFFF));
    // 按位取反操作，由于在SSE指令集中没有需要通过异或实现。

```

最后的时间开销对比，由于在 X86 平台和 ARM 平台下二者硬件参数存在差异，因此比较意义不大，只对串并行算法做对比。

	SIMD	串行
无编译优化	3.40103s	2.64341s
O1	0.859473s	1.1589s
O2	0.869919s	1.09493s
O3	0.630818s	1.07665s

表 4: 不同优化力度下的时间开销

由此可见在 X86 架构下，运行速度显著快于 ARM 架构，一方面是硬件相关参数上的差异，另一方面是二者设计的差异前者为了高性能，而后者则为了低能耗牺牲了部分性能。

在 SSE 下的加速比大约为 **1.70**。低于 ARM 架构。

### 4.2 尝试不同并行度

在原版 SIMD 的基础上，我还尝试了同时处理 2,4,8,16 个指令的 SIMD。出于保证架构一致的考虑，均在 X86 架构下使用 AVX 指令集（这个指令集包含了 SSE）。由于 64 位寄存器现代编译器已经不对其操作进行优化，因此对于两个指令的同时处理，考虑采用 128 位寄存器，只利用其两位的方式实现。对于 8 个以及 16 个指令的 SIMD 则分别使用 AVX 的 256bit 寄存器和 AVX-512 的 512bit 指令集。

二路并行指令集与四路并行的指令集相同，因此不做赘述；

八路并行的 256bit 寄存器的名称为 `__m256i`，十六路并行的 512bit 寄存器名称为 `__m512i`。

实验均在 O3 优化下进行，时间开销以及加速比结果如下：



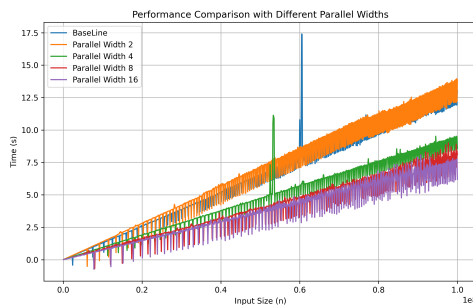
AVX 256-bit (YMM)	AVX-512 512-bit (ZMM)	注释
<code>_mm256_loadu_si256((__m256i*)y);</code>	<code>_mm512_loadu_si512((__m512i*)y);</code>	从内存中加载数据
<code>_mm256_set1_epi32(0x67452301);</code>	<code>_mm512_set1_epi32(0x67452301);</code>	向量广播操作
<code>_mm256_slli_epi32(num,n);</code>	<code>_mm512_slli_epi32(num,n);</code>	左移操作
<code>_mm256_srli_epi32(num,n);</code>	<code>_mm512_srli_epi32(num,n);</code>	右移操作
<code>_mm256_or_si256(x,y);</code>	<code>_mm512_or_si512(x,y);</code>	按位或操作
<code>_mm256_and_si256(x,y);</code>	<code>_mm512_and_si512(x,y);</code>	按位与操作
<code>_mm256_xor_si256(x,y);</code>	<code>_mm512_xor_si512(x,y);</code>	按位异或操作
<code>_mm256_add_epi32(x,y);</code>	<code>_mm512_add_epi32(x,y);</code>	加法
<code>_mm256_xor_si256(x,_mm256_set1_epi32(0xFFFFFFFF));</code>	<code>_mm512_xor_si512(x,_mm512_set1_epi32(0xFFFFFFFF));</code>	按位取反 (通过异或实现)

表 5: AVX/AVX-512 指令对照表 (针对 `uint32_t` 数据类型)

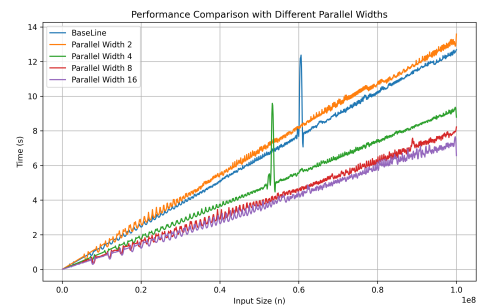
	串行	二路并行	四路并行	八路并行	十六路并行
时间	1.07665s	0.989887s	0.630818s	0.595189s	0.607883s
加速比		1.02	1.70	<b>1.81</b>	1.78

表 6: 不同并行程度下的开销

此外,测试了在不同问题规模下,不同并行路数的时间开销,结果如下:



((a)) 原始不同规模并行度时间开销



((b)) 平滑化后的

图 4.2: 不同规模不同并行度 SIMD 时间开销

我是以 100000 作为初始大小,以步长为 100000 逐步增加到 100000000 从而测试不同并行度算法在不同问题规模下的时间开销。由于受系统波动影响,原始图像呈现出震荡,不便于观察趋势,因此对原始数据的图像做了一定的平滑化操作。

可以发现总体上来说,并行度越高,其性能就越好。

虽然从平滑化后的数据可以发现二路版本的时间开销较串行版本有所增加,但实际上二路并行的版本是优于串行版本的,只是加速比很小。这很大程度上是因为目前支持 64bit(2\*32bit) 的寄存器已经不再被编译优化,因此采用 128bit 的寄存器空出 64bit 的方法实现,但是这样可能会导致内存访问不连续的问题,从而导致其优化性能不佳。

### 4.3 Profiling + 加速成功分析

根据 Table ??, 在无编译优化的情况下, SIMD 版本的时间开销甚至大于串行版本,而在启用编译器优化后 SIMD 的性能则超过了串行版本。

在分析这背后的原理之前,我们先通过下面的表格来了解编译器的 O1,O2,O3 优化。

优化等级	特性说明
<b>O1</b>	启用基础优化，如删除死代码、简化表达式、减少跳转。提高性能同时保证编译速度快，代码稳定。
<b>O2</b>	启用更高级优化，如函数内联、循环优化、跳转预测，更激进的死代码消除。性能与稳定性平衡良好。
<b>O3</b>	启用最激进的优化，包括大规模循环展开、深度内联、自动矢量化（SIMD）、预取指令等，最大化性能但编译时间更长。

表 7: 编译器不同优化等级（O1、O2、O3）特性总结

因为只要启用了编译优化，SIMD 的性能就优于了串行算法，因此我在 Godbolt 网站上生成 SSE 下的 SIMD 算法在无编译优化以及 O1 编译优化下的汇编代码，发现开启 O1 编译优化的代码，相较于没有开启编译优化的，增删改了部分代码。

例如，删去了大量的调用 FF, GG, HH, II 等 hash 操作的汇编指令，相应地将这些步骤直接内联到 MD5\_hash 内部 (红色部分为删除的，黄色的为修改或新增的):

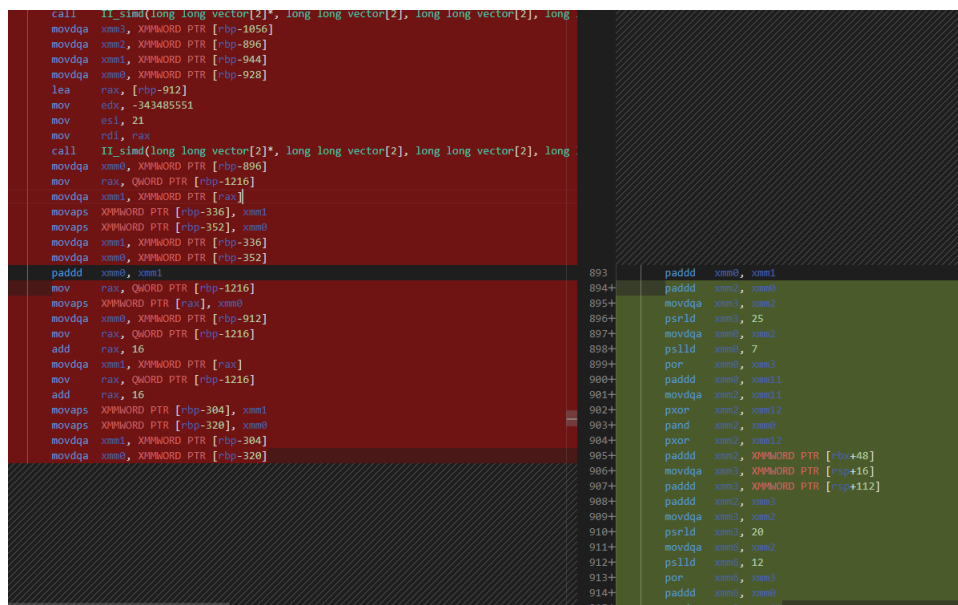


图 4.3: 汇编代码对比

通过函数的内联，减少了调用函数时参数的频繁的出栈入栈，从而大大加速了程序的运行速度，从而实现了加速。

再将 O1, O3 对应的汇编代码作对比，可以明显地看到，O3 采取了更加激进的优化策略，一方面在原来优化策略的基础上将函数完全内联，如 StringProcess 函数被完全内联进了 MD5\_Hash 过程:

```

call    std::_cxx11::basic_string<char, std::char_traits<char>, st
mov     qword ptr [rsp + 16], rax
mov     rcx, qword ptr [rsp + 400]
mov     qword ptr [rsp + 32], rcx
.LBB2_3:
test    rbp, rbp
je      .LBB2_22
cmp     rbp, 1
jne     .LBB2_21
movzx   ecx, byte ptr [r14]
mov     byte ptr [rax], cl
jmp     .LBB2_22
.LBB2_21:
mov     rdi, rax
mov     rsi, r14
mov     rdx, rbp
call    memcpy@PLT
.LBB2_22:
mov     qword ptr [rsp + 80], rbx
mov     rax, qword ptr [rsp + 400]
mov     qword ptr [rsp + 24], rax
mov     rcx, qword ptr [rsp + 16]
mov     byte ptr [rcx + rax], 0
mov     rax, qword ptr [rsp + 16]
mov     qword ptr [rsp], rax
mov     r13, qword ptr [rsp + 24]
lea     eax, [8*r13]
lea     ecx, [8*r13 + 511]
test    eax, eax
cmovns  ecx, eax
and     ecx, -512
sub     eax, ecx
mov     ecx, 448
sub     ecx, eax
cmp     eax, 448
mov     edx, 512

```

(a) O3 优化下函数内联

```

sub     rsp, 424
lea     r14, [rsp+160]
lea     r15, [rsp+176]
.L50:
mov     rsi, QWORD PTR [r13+0+12*8]
mov     rdx, QWORD PTR [r13+8+12*8]
mov     rdi, r14
mov     QWORD PTR [rsp+160], r15
call    void std::_cxx11::basic_string<char, std::char_traits<char>,
lea     rsi, [rsp+96+12]
mov     rdi, r14
call    StringProcess(std::_cxx11::basic_string<char, std::char_traits<
mov     rdi, QWORD PTR [rsp+160]
lea     rbx, [rsp+128]
mov     QWORD PTR [rbx+12*2], rax
cmp     rdi, r15
je      .L48
mov     rax, QWORD PTR [rsp+176]
lea     rsi, [rax+1]
call    operator delete(void*, unsigned long)

```

(b) O1 优化下部分函数未内联

图 4.4: 不同优化下汇编代码差异分析

同时也能看到，在 O3 优化的版本下拥有更多的向量寄存器指令，这表明 O3 优化，会进一步进行循环展开以及尝试 SIMD。这样一方面减少了对内存的访问，同时也提升了运算的效率。

```

pand    xmm1, xmm0
movdqa  xmm1, xmm0
pandn   xmm1, xmm0
por     xmm1, xmm0
movdqa  xmm0, xmmword ptr [rsp + 400]
movdqa  xmmword ptr [rsp + 64], xmm0
movdqa  xmm12, xmmword ptr [rsp + 416]
movdqa  xmm4, xmmword ptr [rsp + 432]
movdqa  xmmword ptr [rsp + 80], xmm0
movdqa  xmmword ptr [rsp], xmm0
movdqa  xmm0, xmmword ptr [rsp + 448]
movdqa  xmmword ptr [rsp + 272], xmm0
movdqa  xmm1, xmm0
padd    xmm1, xmm0
padd    xmm1, xmm0
padd    xmm1, xmm0
padd    xmm1, xmmword ptr [rip + .LCPI2_4]
movdqa  xmm0, xmm1
psrld   xmm0, 25
pslld   xmm1, 7
por     xmm1, xmm0
padd    xmm1, xmm0
movdqa  xmm0, xmm0
pand    xmm0, xmm3
movdqa  xmm1, xmm3
pandn   xmm1, xmm14
por     xmm1, xmm0
movdqa  xmm0, xmmword ptr [rsp]
padd    xmm0, xmm12
padd    xmm0, xmm1
padd    xmm0, xmmword ptr [rip + .LCPI2_5]
movdqa  xmm1, xmm0
psrld   xmm1, 20
pslld   xmm0, 12
por     xmm0, xmm3
padd    xmm0, xmm3
movdqa  xmm1, xmm0
pand    xmm1, xmm3
movdqa  xmm1, xmm14
padd    xmm1, xmm0

```

```

1002+  mov     esi, QWORD PTR [rsp+88]
1003+  movdqa  xmm1, XMMWORD PTR [rsp+160]
1004+  pand    xmm1, xmm0
1005+  pxor    xmm1, XMMWORD PTR [rip+48]
1006+  add     esi, 1
1007+  padd    xmm0, xmm7
1008+  padd    xmm1, xmm13
1009+  padd    xmm1, xmm0
1010+  movdqa  xmm1, xmm1
1011+  pslld   xmm1, 7
1012+  psrld   xmm0, 25
1013+  movdqa  xmm1, xmm1
1014+  movdqa  xmm1, xmm10
1015+  por     xmm0, xmm2
1016+  pxor    xmm1, xmm11
1017+  movdqa  xmm0, XMMWORD PTR [rsp+16]
1018+  padd    xmm1, XMMWORD PTR [rsp+176]
1019+  padd    xmm0, xmm0
1020+  pand    xmm1, xmm0
1021+  pxor    xmm1, xmm11
1022+  padd    xmm1, XMMWORD PTR [rip+48]

```

图 4.5: 更多的向量指令

## 参考文献