

軟體分析與最佳化 Workload Analysis

Stage 2

組別：2

成員：612410017 林靖紳、612410066 蔡宏遠

Review Stage 1 Results

- 根據 Stage 1 的分析結果，我們打算以 gcc 作為之後主要分析的編譯器

- GCC

- -O3

real	0m4.940s
user	0m14.492s
sys	0m3.518s

- -O0

real	0m5.923s
user	0m19.574s
sys	0m5.635s

- ICC

- -O3

real	0m5.249s
user	0m13.975s
sys	0m3.782s

- -O0

real	0m5.344s
user	0m14.020s
sys	0m4.003s

Execution environments

- CPU information

```
ashen@Stephanie-Lin:~$ lscpu
Architecture:          x86_64
CPU op-mode(s):        32-bit, 64-bit
Address sizes:          39 bits physical, 48 bits virtual
Byte Order:             Little Endian
CPU(s):                 12
On-line CPU(s) list:   0-11
Vendor ID:              GenuineIntel
Model name:             11th Gen Intel(R) Core(TM) i5-11500 @ 2.70GHz
CPU family:             6
Model:                 167
Thread(s) per core:    2
Core(s) per socket:    6
Socket(s):              1
Stepping:               1
CPU max MHz:           4600.0000
CPU min MHz:           800.0000
BogoMIPS:               5424.00
```

- Memory

```
ashen@Stephanie-Lin:~$ free -h
```

	total	used	free	shared	buff/cache	available
Mem:	31Gi	4.2Gi	12Gi	1.8Gi	14Gi	24Gi
Swap:	2.0Gi	0B	2.0Gi			

- OS version

```
ashen@Stephanie-Lin:~$ lsb_release -a
```

No LSB modules are available.

Distributor ID: Ubuntu

Description: Ubuntu 22.04.2 LTS

Release: 22.04

Codename: jammy

- GCC version

```
ashen@Stephanie-Lin:~$ gcc --version
```

gcc (Ubuntu 11.4.0-1ubuntu1~22.04) 11.4.0

Copyright (C) 2021 Free Software Foundation, Inc.

This is free software; see the source for copying conditions. There is NO warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.

- ICC version

```
ashen@Stephanie-Lin:~/Documents/Software_Analysis-git/HW2$ icc --version
```

icc: remark #10441: The Intel(R) C++ Compiler Classic (ICC) is deprecated and will be removed from product release in the second half of 2023. The Intel(R) oneAPI DPC++/C++ Compiler (ICX) is the recommended compiler moving forward. Please transition to use this compiler. Use '-diag-disable=10441' to disable this message.

icc (ICC) 2021.10.0 20230609

Copyright (C) 1985-2023 Intel Corporation. All rights reserved.

編譯參數：

```
1 export CFLAGS="-DUNIX -O0 -funroll-loops -fprefetch-loop-arrays ${
```

```
2 export CXXFLAGS="-DUNIX -O0 -funroll-loops -fprefetch-loop-arrays
```

```
3 export CPPFLAGS=""
```

```
4 export CXXCPPFLAGS=""
```

```
5 export LDFLAGS="-L${CC_HOME}/lib64 -L${CC_HOME}/lib"
```

```
6 export LIBS=""
```

```
7 export EXTRA_LIBS=""
```

```
8 export PARMACS_MACRO_FILE="pthread"
```

使用 vtune 初步分析

gcc -O0

Elapsed Time[?]: 5.041s

CPU Time [?] :	25.519s	
Effective Time [?] :	25.519s	
Spin Time [?] :	0s	
Overhead Time [?] :	0s	
Instructions Retired:	170,988,300,000	
Microarchitecture Usage [?] :	51.3%	of Pipeline Slots
CPI Rate [?] :	0.632	
Total Thread Count:	16	
Paused Time [?] :	0s	

Top Hotspots

This section lists the most active functions in your application. Optimizing these hotspot functions typically results in improving overall application performance.

Function	Module	CPU Time [?]	% of CPU Time [?]
[Outside any known module]	[Unknown]	6.112s	24.0%
rabinseg	dedup	4.075s	16.0%
deflate_slow	dedup	3.715s	14.6%
pqdownheap	dedup	3.149s	12.3%
sha1_block_data_order	dedup	2.617s	10.3%
[Others]	N/A*	5.850s	22.9%

*N/A is applied to non-summable metrics.

1. 首先我們使用 gcc -O0 編譯發現 hotspot 為 rabinseg 函式，緊接著是第二高的 hotspot 是 deflate_slow 函式
2. 在這樣的情況下，需要考量優化的函式有兩個，要直接用腦袋想兩個辦法有些麻煩，因此我們先嘗試使用 compiler 的編譯選項做簡單且直接的優化

gcc -O3

Elapsed Time: 4.080s

CPU Time:	16.838s
Effective Time:	16.838s
Spin Time:	0s
Overhead Time:	0s
Instructions Retired:	125,096,400,000
Microarchitecture Usage:	52.6% of Pipeline Slots
CPI Rate:	0.564
Total Thread Count:	16
Paused Time:	0s

Top Hotspots

This section lists the most active functions in your application. Optimizing these hotspot functions typically results in improving overall application performance.

Function	Module	CPU Time	% of CPU Time
deflate_slow	dedup	3.523s	20.9%
pqdownheap	dedup	3.305s	19.6%
[Outside any known module]	[Unknown]	2.660s	15.8%
sha1_block_data_order	dedup	1.878s	11.2%
rabinseg	dedup	1.440s	8.5%
[Others]	N/A*	4.032s	23.9%

*N/A is applied to non-summable metrics.

1. 可以看到 gcc -O3 編譯發現 hotspot 只剩下 deflate_slow，並優化了 rabinseg 函式的執行時間比例
2. 因此後續主要會針對在 -O3 編譯選項下的 deflate_slow 函式進行優化

討論 -O3 如何優化 rabinseg 函式

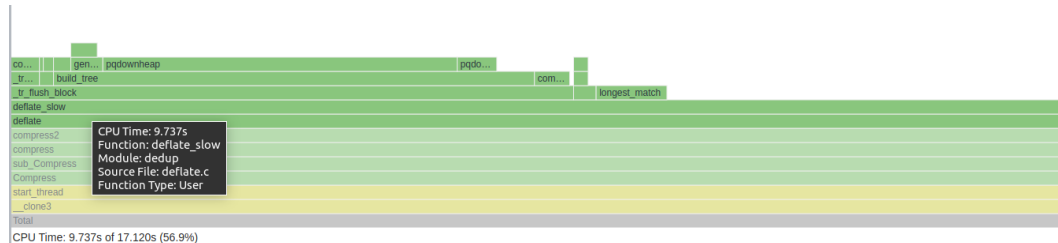
• -O0

72	int rabinseg(uchar *p, int n, int winlen, u32int *rabin)	0s	5.4	0x3844	93	shll \$0x8, -0x10(%rbp)	183.186ms
73	int i;			0x3848	94	movl -0x14(%rbp), %eax	59.735ms
74	u32int h;			0x384b	94	lea 0x1(%rax), %edx	4.978ms
75	u32int x;			0x384e	94	movl %edx, -0x14(%rbp)	188.164ms
76				0x3851	94	movsxd %eax, %rdx	31.858ms
77	USED(winlen);			0x3854	94	movq -0x28(%rbp), %rax	61.726ms
78	if(n < NWINDOW)			0x3858	94	add %rdx, %rax	5.973ms
79	return n;			0x385b	94	movzxb (%rax), %eax	173.230ms
80				0x385e	94	movzx %al, %eax	36.836ms
81	h = 0;			0x3861	94	orl %eax, -0x10(%rbp)	464.934ms
82	for(i=0; i<NWINDOW; i++){	0.006s	29.7	0x3864	95	movl -0xc(%rbp), %eax	102.544ms
83	x = h >> 24;	0.007s	13.5	0x3867	95	lea (,%rax,4), %rdx	1.991ms
84	h = (h<<8) p[i];	0.014s	102.6	0x386f	95	movq -0x38(%rbp), %rax	19.912ms
85	h ^= rabin[i];	0.019s	91.8	0x3873	95	add %rdx, %rax	94.500ms
86	}			0x3876	95	movl (%rax), %eax	140.376ms
87	if((h & RabinMask) == 0)			0x3878	95	xorl %eax, -0x10(%rbp)	697.898ms
88	return i;			0x387b	96	movl -0x10(%rbp), %eax	87.611ms
89	while(i<n){	0.191s	2,157.3	0x387e	96	and \$0xffff, %eax	99.558ms
90	x = p[i-NWINDOW];	0.535s	6,369.3	0x3883	96	test %eax, %eax	0.996ms
91	h ^= rabin[i];	0.391s	3,461.4	0x3885	96	jnz 0x388c <Block 11>	238.938ms
92	x = h >> 24;	0.216s	2,092.5	0x3887		Block 10:	
93	h <<= 8;	0.183s	1,325.7	0x3887	97	movl -0x14(%rbp), %eax	0.996ms
94	h = p[i];	1.027s	5,915.7	0x388a	97	jmp 0x388b <Block 13>	
95	h ^= rabin[i];	1.057s	4,368.6	0x388c		Block 11:	
96	if((h & RabinMask) == 0)	0.427s	2,162.7	0x388c	89	movl -0x14(%rbp), %eax	2.987ms
97	return i;	0.001s		0x388f	89	cmpl -0x2c(%rbp), %eax	
98	}			0x3892	89	jl 0x388b <Block 9>	188.164ms
99	return n;			0x3898		Block 12:	
100	}			0x3898	99	movl -0x2c(%rbp), %eax	
101				0x389b		Block 13:	

• -03

Source Line ▲	Source	CPU Time	Ins	Address ▲	Source Line	Assembly	CPU Tim
59	for(i=0; i<256; i++)			0x39b9	95	mov %eax, %eax	
60	rabinwintab[i] = fpwinreduce(ir			0x39bb	93	shl \$0x8, %ebx	
61	return;			0x39be	95	movl (%rcx,%rax,4), %eax	
62	}			0x39c1	94	or %ebx, %esi	
63				0x39c3	95	xor %esi, %eax	
64	void rabininit(int winlen, u32int * rabin			0x39c5	96	mov \$0x22, %esi	
65	//rabintab = malloc(256*sizeof rabintab			0x39ca	96	test \$0xffff, %eax	
66	//rabinwintab = malloc(256*sizeof rabin			0x39cf	96	jz 0x3b20 <Block 19>	
67	fpmkredtab(irrpoly, 0, rabintab);			0x39d5		Block 8:	
68	fpmkwintredtab(irrpoly, winlen, rabintab			0x39d5	96	cmp \$0x1, %r12	
69	return;			0x39d9	96	jz 0x3a50 <Block 13>	0ms
70	}			0x39db		Block 9:	
71				0x39db	96	cmp \$0x2, %r12	
72	int rabinseg(uchar *p, int n, int winlen,	2.987ms		0x39df	96	jz 0x3a1a <Block 11>	
73	int i;			0x39e1		Block 10:	
74	u32int h;			0x39e1	91	movzxb -0x21(%rdi,%rsi,1), %r9d	
75	u32int x;			0x39e7	94	movzxb -0x1(%rdi,%rsi,1), %ebx	
76				0x39ec	94	mov \$0x23, %esi	
77	USED(winlen);			0x39f1	91	xorl (%r8,%r9,4), %eax	
78	if(n < NWINDOW)			0x39f5	94	mov \$0x22, %r9d	
79	return n;			0x39fb	93	mov %eax, %r12d	
80				0x39fe	92	shr \$0x18, %eax	
81	h = 0;			0x3a01	95	mov %eax, %edx	
82	for(i=0; i<NWINDOW; i++){	0.996ms		0x3a03	93	shl \$0x8, %r12d	
83	x = h >> 24;	0.996ms		0x3a07	95	movl (%rcx,%rdx,4), %eax	
84	h = (h<<8) p[i];	3.982ms		0x3a0a	94	or %r12d, %ebx	
85	h ^= rabintab[x];	3.982ms		0x3a0d	95	xor %ebx, %eax	
86	}			0x3a0f	96	test \$0xffff, %eax	
87	if((h & RabinMask) == 0)	0ms		0x3a14	96	jz 0x3b20 <Block 19>	
88	return i;			0x3a1a		Block 11:	
89	while(i<n){	1.991ms		0x3a1a	91	movzxb -0x21(%rdi,%rsi,1), %r9d	
90	x = p[i-NWINDOW];			0x3a20	94	movzxb -0x1(%rdi,%rsi,1), %ebx	
91	h ^= rabinwintab[x];	61.726ms		0x3a25	91	xorl (%r8,%r9,4), %eax	
92	x = h >> 24;	9.956ms		0x3a29	94	mov %esi, %r9d	
93	h <<= 8;	124.447ms		0x3a2c	94	add \$0x1, %rsi	
94	h = p[i++];	237.942ms		0x3a30	93	mov %eax, %r12d	
95	h ^= rabintab[x];	866.150ms		0x3a33	92	shr \$0x18, %eax	
96	if((h & RabinMask) == 0)	124.447ms		0x3a36	95	mov %eax, %eax	0ms
97	return i;			0x3a38	93	shl \$0x8, %r12d	
98	}			0x3a3c	95	movl (%rcx,%rax,4), %eax	0.996ms
99	return n;			0x3a3f	94	or %r12d, %ebx	
100	}			0x3a42	95	xor %ebx, %eax	
59	for(i=0; i<256; i++)			0x3a6a	93	mov %eax, %r12d	
60	rabinwintab[i] = fpwinreduce(ir			0x3a6d	92	shr \$0x18, %eax	3.982ms
61	return;			0x3a70	93	shl \$0x8, %r12d	41.814ms
62	}			0x3a74	94	or %r12d, %ebx	6.969ms
63				0x3a77	95	xorl (%rcx,%rax,4), %ebx	218.031ms
64	void rabininit(int winlen, u32int * rabin			0x3a7a	95	mov %ebx, %eax	1.991ms
65	//rabintab = malloc(256*sizeof rabintab			0x3a7c	96	lea 0x1(%rsi), %rbx	
66	//rabinwintab = malloc(256*sizeof rabin			0x3a80	96	test \$0xffff, %eax	18.916ms
67	fpmkredtab(irrpoly, 0, rabintab);			0x3a85	96	jz 0x3b20 <Block 19>	
68	fpmkwintredtab(irrpoly, winlen, rabintab			0x3a8b		Block 15:	
69	return;			0x3a8b	91	movzxb -0x21(%rdi,%rbx,1), %r9d	23.894ms
70	}			0x3a91	94	movzxb -0x1(%rdi,%rbx,1), %edx	1.991ms
71				0x3a96	91	xorl (%r8,%r9,4), %eax	1.991ms
72	int rabinseg(uchar *p, int n, int winlen,	2.987ms		0x3a9a	94	mov %ebx, %r9d	12.942ms
73	int i;			0x3a9d	94	lea 0x2(%rsi), %rbx	27.876ms
74	u32int h;			0x3aa1	93	mov %eax, %r12d	1.991ms
75	u32int x;			0x3aa4	92	shr \$0x18, %eax	0ms
76				0x3aa7	93	shl \$0x8, %r12d	23.894ms
77	USED(winlen);			0x3aab	94	or %r12d, %edx	41.814ms
78	if(n < NWINDOW)			0x3aae	95	xorl (%rcx,%rax,4), %edx	174.226ms
79	return n;			0x3ab1	95	mov %edx, %eax	1.991ms
80				0x3ab3	96	test \$0xffff, %eax	28.872ms
81	h = 0;			0x3ab8	96	jz 0x3b20 <Block 19>	0.996ms
82	for(i=0; i<NWINDOW; i++){	0.996ms		0x3aba		Block 16:	
83	x = h >> 24;	0.996ms		0x3aba	91	movzxb -0x1f(%rdi,%rsi,1), %r9d	0ms
84	h = (h<<8) p[i];	3.982ms		0x3ac0	94	movzxb 0x1(%rdi,%rsi,1), %edx	23.894ms
85	h ^= rabintab[x];	3.982ms		0x3ac5	91	xorl (%r8,%r9,4), %eax	0.996ms
86	}			0x3ac9	94	mov %ebx, %r9d	16.925ms
87	if((h & RabinMask) == 0)	0ms		0x3acc	94	lea 0x3(%rsi), %rbx	0.996ms
88	return i;			0x3ad0	93	mov %eax, %r12d	23.894ms
89	while(i<n){	1.991ms		0x3ad3	92	shr \$0x18, %eax	2.987ms
90	x = p[i-NWINDOW];			0x3ad6	93	shl \$0x8, %r12d	21.903ms
91	h ^= rabinwintab[x];	61.726ms		0x3ada	94	or %r12d, %edx	0ms
92	x = h >> 24;	9.956ms		0x3add	95	xorl (%rcx,%rax,4), %edx	238.938ms
93	h <<= 8;	124.447ms		0x3ae0	95	mov %edx, %eax	0.996ms
94	h = p[i++];	237.942ms		0x3ae2	96	test \$0xffff, %eax	30.863ms
95	h ^= rabintab[x];	866.150ms		0x3ae7	96	jz 0x3b20 <Block 19>	
96	if((h & RabinMask) == 0)	124.447ms		0x3ae9		Block 17:	
97	return i;			0x3ae9	91	movzxb -0x1e(%rdi,%rsi,1), %r9d	
98	}			0x3aef	91	xorl (%r8,%r9,4), %eax	16.925ms
99	return n;			0x3af3	94	mov %ebx, %r9d	
100	}			0x3af6	91	mov %eax, %edx	17.920ms

- 從 flame graph 中可以看到以下幾點：
 - 整個程式的執行流程
 - 函數呼叫的頻率與深度
 - 深度較深的函數呼叫可能是潛在的性能瓶頸，像是 deflate_slow 函式
 - 耗時的操作
 - 矩形區塊的寬度反映了相應函數執行的時間



Sampling

CPI Rate

可以從 CPI 看到 -O3 程式的平行化程度有明顯的提昇

- CPI (-O0) : 0.632
- CPI (-O3) : 0.564

Memory

可以看到下圖使用 -O3 與 -O0 比較在 load/store 還有 LLC miss 的次數都大幅度的降低

- gcc -O0

Elapsed Time	4.792s	
CPU Time	20.030s	
Memory Bound	7.1%	of Pipeline Slots
L1 Bound	14.1%	of Clockticks
L2 Bound	2.0%	of Clockticks
L3 Bound	1.9%	of Clockticks
DRAM Bound	1.0%	of Clockticks
DRAM Bandwidth Bound	0.0%	of Elapsed Time
Store Bound	1.4%	of Clockticks
Loads	53,766,070,468	
Stores	24,477,828,827	
LLC Miss Count	394,603	
Total Thread Count	16	
Paused Time	0s	

- gcc -O3

Elapsed Time: 3.866s

CPU Time:

11.712s

Memory Bound:

8.3%

of Pipeline Slots

L1 Bound:

16.0%

of Clockticks

L2 Bound:

0.0%

of Clockticks

L3 Bound:

0.9%

of Clockticks

DRAM Bound:

1.4%

of Clockticks

DRAM Bandwidth Bound:

0.0%

of Elapsed Time

Store Bound:

2.3%

of Clockticks

Loads:

20,190,359,746

Stores:

8,824,850,780

LLC Miss Count:

329,261

Total Thread Count:

16

Paused Time:

0s

是否有使用 compiler directive 的機會

在 deflate_slow 函式中使用率最高的 hotspot 就是 INSERT_STRING，這個函式主要是將 str 插入到 hash dictionary 中，並記錄具有相同 hash key 的最近字串的資訊

Source Line	Source	CPU Time: Total	CPU	Address	Source Line	Assembly	CPU Time: Total	CPU
1553	/*			0xd544	1557	pushq %r12		
1554	local block_state deflate_slow(s, flush)			0xd546	1623	lea 0x120d3(%rip), %r12		
1555	deflate_state *s;			0xd54d	1557	pushq %rbp		
1556	int flush;			0xd54e	1557	mov %esi, %ebp		
1557	{			0xd550	1557	pushq %rbx		
1558	Ipos hash_head = NIL; /* head of hash chain */			0xd551	1557	mov %rdi, %rbx		
1559	int bflush; /* set if current block is			0xd554	1557	sub \$0x18, %rsp		
1560				0xd558	1568	movl 0xa4(%rdi), %esi		
1561	/* Process the input block. */			0xd55e	1568	data16 nop		
1562	for (;;) {			0xd560		Block 2:		
1563	/* Make sure that we always have enough lookahead			0xd560	1568	cmp \$0x105, %esi		
1564	* at the end of the input file. We need MAX_MATCH			0xd566	1568	jbe 0xd70e <Block 22>		
1565	* for the next match, plus MIN_MATCH bytes to			0xd56c		Block 3:		
1566	* string following the next match.			0xd56c	1580	movl 0x9c(%rbx), %eax	0.2%	
1567	*/			0xd572		Block 4:		
1568	if (s->lookahead < MIN_LOOKAHEAD) {			0xd572	1580	movq 0x50(%rbx), %rcx	0.5%	
1569	fill_window(s);	0.1%		0xd576	1580	lea 0x2(%rax), %esi	0.1%	
1570	if (s->lookahead < MIN_LOOKAHEAD && flush			0xd579	1580	movl 0x70(%rbx), %edi	0.1%	
1571	return need_more;			0xd57c	1580	mov %eax, %r11d		
1572	}			0xd57f	1580	movq 0x60(%rbx), %r9		
1573	if (s->lookahead == 0) break; /* flush the			0xd583	1580	movq 0x60(%rbx), %r15		
1574	}			0xd587	1580	movzxb (%rcx,%rsi,1), %r8d		
1575				0xd58c	1580	movl 0x80(%rbx), %ecx	0.5%	
1576	/* Insert the string window[strstart .. strstart			0xd592	1580	andl 0x4c(%rbx), %r11d	0.4%	
1577	* dictionary, and set hash_head to the head of			0xd596	1580	shl %cl, %edi	0.1%	
1578	*/			0xd598	1580	xor %edi, %r8d	0.2%	
1579	if (s->lookahead >= MIN_MATCH) {			0xd59b	1580	andl 0x7c(%rbx), %r8d		
1580	INSERT_STRING(s, s->strstart, hash_head);	11.6%		0xd59f	1580	lea (%r9,%r8,2), %r10	0.1%	
1581	}			0xd5a3	1580	movl %r8d, 0x70(%rbx)	0.4%	
1582				0xd5a7	1580	movzxb (%r10), %r14d	0.5%	
1583	/* Find the longest match, discarding those <			0xd5ab	1580	movw %r14w, (%r15,%r11,2)	6.4%	
1584	*/			0xd5b0	1580	movw %ax, (%r10)	2.0%	
1585	s->prev_length = s->match_length, s->prev_match	0.4%		0xd5b4		Block 5:		
1586	s->match_length = MIN_MATCH-1;			0xd5b4	1585	movl 0x90(%rbx), %edx	0.4%	

且在 INSERT_STRING 函式的部份已使用 compiler directive

```
/* =====
 * Insert string str in the dictionary and set match_head to the previous head
 * of the hash chain (the most recent string with same hash key). Return
 * the previous length of the hash chain.
 * If this file is compiled with -DFASTEST, the compression level is forced
 * to 1, and no hash chains are maintained.
 * IN assertion: all calls to to INSERT_STRING are made with consecutive
 * input characters and the first MIN_MATCH bytes of str are valid
 * (except for the last MIN_MATCH-1 bytes of the input file).
 */
#ifdef FASTEST
#define INSERT_STRING(s, str, match_head) \
    (UPDATE_HASH(s, s->ins_h, s->window[(str) + (MIN_MATCH-1)]), \
     match_head = s->head[s->ins_h], \
     s->head[s->ins_h] = (Pos)(str))
#else
#define INSERT_STRING(s, str, match_head) \
    (UPDATE_HASH(s, s->ins_h, s->window[(str) + (MIN_MATCH-1)]), \
     match_head = s->prev[(str) & s->w_mask] = s->head[s->ins_h], \
     s->head[s->ins_h] = (Pos)(str))
#endif
```

是否有使用 SIMD 的機會

同樣是 INSERT_STRING 作為 hotspot，反組譯後在 vtune 上並沒有看到使用 SIMD 的痕跡

Source Line ▲	Source	🔥 CPU Time: Total CPU	Address ▲	Source Line	Assembly	🔥 CPU Time: Total CPU
1553	/* Process the input block. */		0xd544	1557	pushq %r12	
1554	local block_state deflate_slow(s, flush)		0xd546	1623	lea 0x120d3(%rip), %r12	
1555	deflate_state *s;		0xd54d	1557	pushq %rbp	
1556	int flush;		0xd54e	1557	mov %esi, %ebp	
1557	{		0xd550	1557	pushq %rbx	
1558	IPos hash_head = NIL; /* head of hash chain */		0xd551	1557	mov %rdi, %rbx	
1559	int bflush; /* set if current block is		0xd554	1557	sub \$0x18, %rsp	
1560			0xd558	1568	movl 0xa4(%rdi), %esi	
1561	/* Process the input block. */		0xd55e	1568	data16 nop	
1562	for (;;) {		0xd560		Block 2:	
1563	/* Make sure that we always have enough lookahead		0xd560	1568	cmp \$0x105, %esi	
1564	* at the end of the input file. We need MAX_MATCH		0xd566	1568	jbe 0xd70e <Block 22>	
1565	* for the next match, plus MIN_MATCH bytes to		0xd56c		Block 3:	
1566	* string following the next match.		0xd56c	1580	movl 0x9c(%rbx), %eax	0.2%
1567	*/		0xd572		Block 4:	
1568	if (s->lookahead < MIN_LOOKAHEAD) {		0xd572	1580	movq 0x50(%rbx), %rcx	0.5%
1569	fill_window(s);	0.1%	0xd576	1580	lea 0x2(%rax), %esi	0.1%
1570	if (s->lookahead < MIN_LOOKAHEAD && flush		0xd579	1580	movl 0x70(%rbx), %edi	0.1%
1571	return need_more;		0xd57c	1580	mov %eax, %r11d	
1572	}		0xd57f	1580	movq 0x60(%rbx), %r9	
1573	if (s->lookahead == 0) break; /* flush the		0xd583	1580	movq 0x60(%rbx), %r15	
1574	}		0xd587	1580	movzxb (%rcx,%rsi,1), %r8d	
1575			0xd58c	1580	movl 0x80(%rbx), %ecx	0.5%
1576	/* Insert the string window[strstart .. strstart		0xd592	1580	andl 0x4c(%rbx), %r11d	0.4%
1577	* dictionary, and set hash_head to the head of		0xd596	1580	shl %cl, %edi	0.1%
1578	*/		0xd598	1580	xor %edi, %r8d	0.2%
1579	if (s->lookahead >= MIN_MATCH) {		0xd59b	1580	andl 0x7c(%rbx), %r8d	
1580	INSERT_STRING(s, s->strstart, hash_head);	11.6%	0xd59f	1580	lea (%r9,%r8,2), %r10	0.1%
1581	}		0xd5a3	1580	movl %r8d, 0x70(%rbx)	0.4%
1582			0xd5a7	1580	movzxb (%r10), %r14d	0.5%
1583	/* Find the longest match, discarding those <		0xd5ab	1580	movw %r14w, (%r15,%r11,2)	6.4%
1584	*/		0xd5b0	1580	movw %ax, (%r10)	2.0%
1585	s->prev_length = s->match_length, s->prev_match	0.4%	0xd5b4		Block 5:	
1586	s->match_length = MIN_MATCH-1;		0xd5b4	1585	movl 0x90(%rbx), %edx	0.4%

但我認為使用 SIMD 的機會是有的，可以將需要插入 hash dictionary 的 string 和 string 準備插入的相關資訊 buffer 起來，然後在批次記錄這些具有相同 hash key 的字串，因此就可以使用 SIMD 指令做插入，大幅度提高 INSERT_STRING 的並行處理效率

討論

- 結論：以上分析出 gcc -O3 中 deflate_slow 函式是整個執行過程中的熱點
- 觀察 deflate_slow 函式的 source code

```
1554 local block_state deflate_slow(s, flush)
1555     deflate_state *s;
1556     int flush;
1557 {
1558     IPos hash_head = NIL;    /* head of hash chain */
1559     int bflush;              /* set if current block must be flushed */
1560
1561     /* Process the input block. */
1562     for (;;) {
1563         /* Make sure that we always have enough lookahead, except
1564          * at the end of the input file. We need MAX_MATCH bytes
1565          * for the next match, plus MIN_MATCH bytes to insert the
1566          * string following the next match.
1567          */
1568         if (s->lookahead < MIN_LOOKAHEAD) {
1569             fill_window(s);
1570             if (s->lookahead < MIN_LOOKAHEAD && flush == Z_NO_FLUSH) {
1571                 return need_more;
1572             }
1573             if (s->lookahead == 0) break; /* flush the current block */
1574         }
1575
1576         /* Insert the string window[strstart .. strstart+2] in the
1577          * dictionary, and set hash_head to the head of the hash chain:
1578          */
1579         if (s->lookahead >= MIN_MATCH) {
1580             INSERT_STRING(s, s->strstart, hash_head);
1581         }
1582     }
```

- 目的：用於實現「slow」壓縮策略
- 大致可以分成以下幾個步驟
 - 透過迴圈處理輸入的 blocks
 - 確保始終有足夠的 lookahead 進行處理，需要時填充window
 - 插入當前字串到字典中
 - 尋找最長匹配，捨棄比前一個匹配短的匹配
 - 根據壓縮策略處理匹配和文字，更新 hash table
 - 在需要時刷新 block

認為可以優化的地方

- 如果是對於包含較多重複資料（字串）的數據，我們認為或許不需要用到 hash table 相對較為複雜的資料結構
- 只需要針對 Run-Length Encoding，簡單地比較相鄰的字節，而不需要複雜的匹配條件。
- 利用相同字節的重複次數即為匹配的長度，不需要進行額外的計算。
- 不用像 deflate_slow 一樣嘗試找到最長的匹配，並僅在確定在下一個窗口位置沒有更好匹配時，才使用當前的匹配
- 影響
 - 壓縮率：
 - 以上的改動可能會影響資料的壓縮效率
 - 以上方法應該比較適用於像是對於像素圖像中的單色區域，或者包含大量重複字符的文本
 - 複雜度：
 - 上述方法的實現應該較為簡單易懂，deflate_slow 有太多的條件判斷分支
- Pseudo code

```
1  deflate_rle:
2      while not end_of_input:
3          fill_window(s)
4          while s.lookahead >= MAX_MATCH:
5              run = find_consecutive_repeats(s)
6              if run >= MIN_MATCH:
7                  process_consecutive_repeats(s, run)
8              else:
9                  process_single_literal(s)
10             if flushing_needed:
11                 flush_block(s)
12
13         flush_block(s, Z_FINISH)
14         return Z_FINISH ? finish_done : block_done
15
```