Optimizing zlib for orm A deflated story

Adenilson Cavalcanti BS. MSc. Staff Engineer - **Arm** San Jose (CA)

THE CONFERENCE MORNING SESSION



WWW. PHDCOMICS. COM

JORGE CHAM @ 201

What to optimize in Chromium



What to optimize in Chromium

- Too big.
- Too many areas.
- What would be helpful?



What to optimize in Chromium

Bulk of content still is:

- Text.
- Images.







PNG

- Powerful format: Palette, pre-filters, compressed.
- Encoder affects behavior.
- Libpng and zlib are 'Bros!'.

Meet Mr. Parrot



Source: <u>https://upload.wikimedia.org/wikipedia/commons/3/3f/ZebraHighRes.png</u>

Parrots are not created equal





Features affect hotspots

== Imag	je has pre-co	mpression filt	ters (2.7MB) ==	
Lib	Command	Shared0bj	method CP	U (%)
zlib	TileWorker	liblink	inflate_fast	1.96
zlib	TileWorker	libblnk	adler32	0.88
blink	TileWorker	liblink	<pre>ImageFrame::setRGBAPremultiply</pre>	0.45
blink	TileWorker	liblink	<pre>png_read_filter_row_up</pre>	0.03*
== Imag	ge was optimi	zed using zop	fli (2.6MB) ==	
Lib	Command	SharedObj	method CP	U (%)
zlib	TileWorker	liblink	inflate fast	3.06
zlib	TileWorker	libblnk	adler32	1.36
blink	TileWorker	liblink	<pre>ImageFrame::setRGBAPremultiply</pre>	0.70
blink	TileWorker	liblink	<pre>png_read_filter_row_up</pre>	0.48*

== Imag	e has no pre	-compression	filters (0.9MB) ==	
Lib	Command	Shared0bj	method CF	U (%)
libpng	TileWorker	liblink	cr png do expand palette	0.88
zlib	TileWorker	liblink	inflate_fast	0.62
blink	TileWorker	liblink	<pre>ImageFrame::setRGBAPremultiply</pre>	0.49
zlib	TileWorker	libblnk	adler32	0.31

NEON: Advanced SIMD (Single Instruction Multiple Data)

SIMD Architecture



- **Optional** in Armv7.
- Mandatory in Armv8.

Registers@Armv7

- 16 registers@128 bits: Q0 Q15.
- 32 registers@64bits: D0 D31.
- Varied set of instructions: load, store, add, mul, etc.

Registers@Armv8 (SIMD&FP, V0 - V31)

- 32 registers@128 bits: Q0 V31.
- 32 registers@64bits: D0 D31.
- 32 registers@32bits: S0 S31.
- 32 registers@8bits: H0 H31.
- Varied set of instructions: load, store, add, mul, etc.

An example: VADD.I16 Q0, Q1, Q2



Candidates

- Inflate_fast: **zlib**.
- Adler32: zlib.
- ImageFrame: Blink.
- png_do_expand_palette: libpng.

Why zlib?

Zlib

Used everywhere (libpng, Skia, freetype, **cronet**, blink, chrome, linux kernel, etc).

Old code base released in 1995.

Written in K&R C style.

Context

Lacks any optimizations for ARM CPUs.

Problem statement

Identify potential optimization candidates and verify positive effects in Chromium.

Potential problems

- Viability of optimization.
- Positive effects.
- Upstreaming.

Implementation

Adler-32

 $A = 1 + D_1 + D_2 + \ldots + D_n \pmod{65521}$ $B = (1 + D_1) + (1 + D_1 + D_2) + \ldots + (1 + D_1 + D_2 + \ldots + D_n)$ (mod 65521) $= n \times D_1 + (n-1) \times D_2 + (n-2) \times D_3 + \ldots + D_n + n \pmod{65521}$

 $Adler-32(D) = B \times 65536 + A$

Adler-32: simplistic implementation

```
// From: https://en.wikipedia.org/wiki/Adler-32
const int MOD_ADLER = 65521;
unsigned long naive_adler32(unsigned char *data,
                             unsigned long len)
    uint32_t a = 1, b = 0;
    unsigned long index;
    for (index = 0; index < len; ++index) {</pre>
        a = (a + data[index]) % MOD_ADLER;
        b = (b + a) \% MOD ADLER;
    return (b << 16) | a;
```

Problems

- Zlib's Adler-32 was more than **7x faster** than naive implementation.
- It is hard to vectorize the following computation:
 void accum(uint32_t *pair, const unsigned char *buf, unsigned int len)

```
unsigned int i;
for (i = 0; i < len; ++i) {
    pair[0] += buf[i];
    pair[1] += pair[0];[
}
```

Problems: how to represent pair[1] or 'B'?

```
unsigned int i;
for (i = 0; i < len; ++i) {
    pair[0] += buf[i];
    pair[1] += pair[0];[
}</pre>
```

Problems: how to represent pair[1] or 'B'?

$$A = 1 + D_1 + D_2 + \ldots + D_n \pmod{65521}$$

$$B = (1 + D_1) + (1 + D_1 + D_2) + \ldots + (1 + D_1 + D_2 + \ldots + D_n)$$

(mod 65521)

$$= n \times D_1 + (n-1) \times D_2 + (n-2) \times D_3 + \ldots + D_n + n \pmod{65521}$$

 $Adler - 32(D) = B \times 65536 + A$

Highly technical drawing (Jan 2017)



Highly technical drawing (Jan 2017)



'Taps' to the rescue

Assembly: https://godbolt.org/g/KMeBAJ

```
static const uint8_t taps[32] = {
    32, 31, 30, 29, 28, 27, 26, 25,
    24, 23, 22, 21, 20, 19, 18, 17,
    16, 15, 14, 13, 12, 11, 10, 9,
    8, 7, 6, 5, 4, 3, 2, 1_};
```

```
uint32x2_t adacc2, s2acc2, as;
uint8x16_t t0 = vld1q_u8(taps), t1 = vld1q_u8(taps + 16);
```

```
uint32x4_t adacc = vdupq_n_u32(0), s2acc = vdupq_n_u32(0);
adacc = vsetq_lane_u32(s[0], adacc, 0);
s2acc = vsetq_lane_u32(s[1], s2acc, 0);
```

Happy end! Up to 18% performance gain in PNG

commit d400b38e450a71e9ed75bd4c1dda6329267e9ce0
Author: Adenilson Cavalcanti <adenilson.cavalcanti@arm.com>
Date: Mon Jan 30 15:30:38 2017 -0800

NEON implementation for Adler32

The checksum is calculated in the uncompressed PNG data and can be made much faster by using SIMD.

Tests in ARMv8 yielded an improvement of about 3x (e.g. walltime was 350ms x 125ms for a 4096x4096 bytes executed 30 times).

This alone yields a performance boost for PNG decoding ranging from 5% to 18% depending on a few factors (SoC, battery status, big/little, etc).

https://bugs.chromium.org/p/chromium/issues/detail?id=688601

Inffast (Simon Hosie)

- Second candidate in the perf profiling was inflate_fast.
- Very high level idea: perform long loads/stores in the byte array.
- Major gains: up to 30% faster!



```
*/
        out = chunkcopy_safe(out, from, len, limit);
else {
    from = out - dist;
                                /* copy direct from output */
                                /* minimum length is three */
    do {
        *out++ = *from++;
        *out++ = *from++;
        *out++ = *from++;
        len -= 3;
    } while (len > 2);
    if (len) {
        *out++ = *from++;
        if (len > 1)
            *out++ = *from++;
    /* Whole reference is in range of current output. No
       range checks are necessary because we start with room
       for at least 258 bytes of output, so unroll and roundoff
       operations can write beyond `out+len` so long as they
       stay within 258 bytes of `out`.
    out = chunkcopy_lapped_relaxed(out, dist, len);
```

Libpng (Richard Townsend)

- NEON optimization in libpng.
- From 10 to 30% improvement.
- Depends on png using a palette.

Impact

Combined effect of 3 patches



Chrome trace: vanilla Nexus6@2014 (116ms)



Chrome trace: patched (73ms) 1.6x improvement



Comparing Arm x Intel



Source: https://commons.wikimedia.org/wiki/File:Apple_and_Orange_-_they_do_not_compare.jpg

Keeping in mind

- Snapdragon[™] 805 @2014.
- 2.7Ghz Krait[™] 450.
- 2MB L2 cache
- 28nm lithography.
- Cellphone.
- EAS kernel.

- 5Y10C launched @2015.
- 2Ghz Intel m5.
- 4MB cache.
- 14nm lithography.
- Ultrabook.
- Regular linux kernel.

Chrome trace: Intel m5@2016 (66ms)



ImageFrameGenerator::decode

Effect of NEON optimization in Zlib



Lessons learned

- **arm** cores can benefit **a lot** from NEON optimizations.
- Performance gains of 2 generations of silicon.
- It pays off to work in a lower software layer (e.g. zlib/libpng).

Happy end? Not yet...

- Requested to perform a study comparing zlibs forks.
- Upstream ARM optimizations.
- Move Chromium to a new/better maintained zlib.

Happy end? Not yet...

- Requested to perform a study comparing zlibs forks. Done!
 - o <u>https://goo.gl/ZUoy96</u>
- Upstream ARM optimizations. Done!
 - <u>https://github.com/Dead2/zlib-ng/commit/ec02ecf104e1d3f183</u>
 <u>6a908a359f20aa93494df5</u>
- Move Chromium to a new/actively maintained zlib.
 - Upgraded/moved PDFium to Chromium's zlib.
 - Zlib-ng didn't release a stable release.



Change of strategy

Sea	rch for messages	- Q		0		
+	POST REPLY C	<u>•</u> • •	۵	*		
link-d et's	ev > optimize zlib for ARM s by 3 authors 🐨 💽					
:	Adenilson Cavalcanti	Aug 14	*	•		
*	Zlib is a compression library used by Chromium code base and its dependencies (skia, libpng, pdfium, freetype, etc) for quite a few tasks ranging from image handling, loading extensions and accessing compressed content (i.e. Content-Encoding: gzip).					
	It is an impressive feat of engineering and it is used all over the place (e.g. the need to support long gone compil	g considering that it is 22 linux kernel). Due to its h lers and operating syster	years o istory ai ns, its	ld nd		

NEON inffast: featured in M62



https://bugs.chromium.org/p/chromium/issues/detail?id=697280

cronet: NEON != ARMv6

	LAIES1: 10.17	UPDAIE)
	CHANGES IN VER THE CPU NO LONG WHEN YOU HOLD D	ESION 10.17: DER OVERHEATS DOWN SPACEBAR.
	COMME	DNT5:
	LONGTIME USERY W	RITES:
	THIS UPDATE BROKE MY CONTROL KEY IS 50 I HOLD SPACEBA CONFIGURED EMACS RAPID TEMPERATURE	E MY WORKFLOW! 5 HARD TO REACH, R INSTEAD, AND I TO INTERPRET A RISE AS CONTROL.
	ADMIN WRITES: THAT'S HORRIFYIN	IG.
	LOOK, MY SETUP W JUST ADD AN OPTION SPACEBAR HEATING.	ites: 10rks for Me. N To REENABLE

EVERY CHANGE BREAKS SOMEONE'S WORKFLOW.

After re-landing... An internal app was broken.



EVERY CHANGE BREAKS SOMEONE'S WORKFLOW.

Second revert (i.e. revert-revert-revert)



Misha Efimov@Google found the bug in the Java app client last Wednesday (Sep 27th).

Re-re-landed on Thur 28th



re-land

What comes next

- Land Adler-32 optimization* (Noel Gordon@Google implemented the same algorithm for Intel).
- Land the libpng optimization.
- CRC32: Armv8 instruction is about 10x faster.
- Compression comes next.

Adler-32 landed on Fri 29th



Adler-32

What comes next

- Land the libpng optimization.
- CRC32: ARMv8 instruction is about 10x faster.
- Fix infback corner case.
- Compression comes next.

Zlib users should consider migrating to Chromium's zlib.

Special Thanks

- Igalia for the invite (Xabier Rodriguez Calvar).
- Arm for sponsoring the trip.
- Chris Blume@Google.
- Team **Arm@UK**: Dave Rodgman, Matteo Franchin, Richard Townsend, Stephen Kyle.
- Team **Arm@US**: Amaury Leleyzour, Simon Hosie.
- Compiler explorer: <u>https://godbolt.org</u>

Questions?



The Arm Trademarks featured in this presentation are registered trademarks or

trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights

reserved. All other marks featured may be trademarks of their respective owners

https://www.arm.com/company/policies/trademarks