SIMD-enhanced libc string functions how it's done

Robert Clausecker <fuz@FreeBSD.org>

Common Tasks

- copying strings (*strcpy*, *memcpy*, ...)
- finding string length (strlen, strnlen, ...)
- finding characters (strchr, memchr, ...)
- comparing strings (strcmp, memcmp, ...)
- finding substrings (*strstr*, *memmem*, ...)
- splitting at delimiters (strspn, strcspn, ...)

Common Tasks

- copying strings (*read* then *write*)
- finding string length (read then compare)
- finding characters (read then compare)
- comparing strings (read then compare)
- finding substrings (complicated)
- splitting at delimiters (read then set match)

What does that mean?

read

- char by char until end of string
- one load/compare/conditional branch per character

write

- char by char until end of string
- one write per character

compare

- char by char until match or end of string
- one compare/conditional branch per character

What does that mean?

read

- char by char until end of string
- one load/compare/conditional branch per character (slow)
 write
 - char by char until end of string
 - one write per character (**slow**)

compare

- char by char until match or end of string
- one compare/conditional branch per character (**slow**)

Conclusion

Conclusion

Strings suck

What can we do about that?

- Get rid of strings (oof...)
- special-purpose instructions (arch dependent)
 - speed varies depending on CPU model
 - often only memcpy(), memset() supported
- strange hacks (hmm...)

SIMD

Your new best friend

SIMD

- Single Instruction Multiple Data
- *SIMD register*: short arrays of numbers
- common lengths: 16, 32, 64 bytes
- same operation on all elements
- but as fast as scalar operations
- SIMD with 16 bytes: 16x scalar performance

Scalar vs. SIMD



Scalar vs. SIMD



typical SIMD operations

Arithmetic (integer/FP)

• addition, subtraction, multiplication, ...

Logic

• element-wise comparison, and, or, xor, ...

Data transfer

• read, write, extract masks, ...

... many more

Strings and SIMD

How can this help us with string processing?

Strings and SIMD

How can this help us with string processing?

- load multiple characters at once
- process them simultaneously
- •
- profit?

f	0	0	b	а	r	b	а	z	\0	

f	0	0	b	а	r	b	а	Z	\0	

f	0	0	b	а	r	b	а	Z	\0	

f	0	0	b	а	r	b	а	Z	\0	

- We can easily overshoot the string's end
- For nul-terminated strings, we won't know where that is until we see the nul byte
- Do we have to iterate char-by-char after all?

What can we do?

What can we do?

String bounds are fictious

What can we do?

String bounds are fictious Let's overcome them!

Overcoming array bounds

- the computer does not know what an array is
- it only knows that there's memory at some addresses but not at others

Overcoming array bounds

- the computer does not know what an array is
- it only knows that there's memory at some addresses but not at others

thus:

- if we don't go too far out of bounds, it'll be fine!
- C doesn't let us, so let's use assembly

How far is too far?

- Memory is organised in *pages*
- size: arch dependent, usually 4096 bytes
- pages are either accessible entirely or not at all
- there is no more fine-grained memory protection
 - (check out CHERI, it's cool)

How far is too far?

- if at least one byte of a string is on a page, the whole page is accessible
- aligned accesses never cross page boundaries

thus:

• if we're careful, it might just work!



	-					-	-				
	f	ο	ο	b	а	r	b	а	z	\0	







Can't use the same approach:

• overreads are fine, overwrites are no good

Can't use the same approach:

• overreads are fine, overwrites are no good

Instead

- write (possibly unaligned) chunks
- last write may overlap previous writes









Can't use the same approach

- strings may have different misalignment
- can't fix this after loading with SSE2

Can't use the same approach

- strings may have different misalignment
- can't fix this after loading with SSE2

Instead

- do aligned reads to check for nul bytes
- then unaligned reads to compare characters



	f	о	ο	b	а	r	b	а	z	\0		
		-			-	-				-	-	

























1		1		I Contraction of the second
	1			
	I			





1				
I				
		•		



f o \0				
		f	ο	\0



copy to bounce buffer on stack





Set Matching

strcspn("foo bar", " \t\n");

- matches each char in string against set
- portable approach: Muła / Langdale algorithm http://0x80.pl/articles/simd-byte-lookup.html
- can we do better?

Set Matching

The Intel way: **pcmpistrm**

- packed compare implicity-terminated string, return mask
- set matching and lots of other features
- conveniently also checks for nul terminators
- probably also brews coffee if you ask nicely



Substring Matching

- That means *strstr(*), *memmem(*)
- really tricky
- most fancy algorithms are optimised for long strings, but our strings are usually short
- wip

Results (amd64)



GB/S

AArch64 port

- work by Getz Mikalsen
- complements the ARM Optimized Routines
- some functions redone due to poor perf
- others not present in ARM's codebase
- assembly functions: strcmp(), strspn(), strcspn(), strlcpy(), strncmp(), memccpy(), strlen(), timingsafe_bcmp(), timingsafe_memcp()
- new wrappers: strpbrk(), strsep(), strcat(), strncat(), strlcat(), bcopy(), bzero()

AArch64 difficulties

- no easy way to get syndrome bitmask from syndrome vector
 - amd64: pmovmskb
 - aarch64: *shrn* + *fmov* to get nibble mask
 - other tricks when only zero / nonzero is needed
- no *pcmpistrm*, Muła/Langdale algorithm instead
- can't count trailing zeroes, only leading zeroes

Current Progress

- 2023 rework of the libc string functions for amd64
 - paid by The FreeBSD Foundation
 - almost all of <string.h>
 - for amd64 baseline (SSE2), some for x86-64-v2
 - landed for 14.1-RELEASE
- later ports as part of GSoC 2024
 - AArch64 by getz@ (landed for upcoming 15-RELEASE)
 - riscv64 by strajabot@ (under acceptance testing)
- Future work: *strstr*(), CHERI / Morello support, PPC64?