Translating C To Rust: Lessons from a User Study

<u>Ruishi Li*</u>, Bo Wang*, Tianyu Li , Prateek Saxena, Ashish Kundu[†]

National University of Singapore Cisco Research[†]



*Equal Contribution

A "dream problem": Translating C to Rust for memory safety



[1] "Translate all C to Rust" by DARPA. USA. 2024

[2] Rust for Linux Github organization: https://github.com/Rust-for-Linux by Linux. 2024

^[3] Rust for Windows: https://github.com/microsoft/windows-rs by Microsoft. 2024

^[4] Bare-metal Rust in Android: <u>https://security.googleblog.com/search/label/rust</u> by Google. 2024

A "dream problem": Translating C to Rust for memory safety



[1] "Translate all C to Rust" by DARPA. USA. 2024

[2] Rust for Linux Github organization: https://github.com/Rust-for-Linux by Linux. 2024

^[3] Rust for Windows: https://github.com/microsoft/windows-rs by Microsoft. 2024

^[4] Bare-metal Rust in Android: <u>https://security.googleblog.com/search/label/rust</u> by Google. 2024

A "dream problem": Translating C to Rust for memory safety



[1] "Translate all C to Rust" by DARPA. USA. 2024

[2] Rust for Linux Github organization: https://github.com/Rust-for-Linux by Linux. 2024

^[3] Rust for Windows: https://github.com/microsoft/windows-rs by Microsoft. 2024

^[4] Bare-metal Rust in Android: https://security.googleblog.com/search/label/rust by Google. 2024

Limitations of existing work

Compiler-based translators (Laertes^[1], Crown^[2])

LLMs-based translators (Flourine^[3], Vert^[4])

[1] Emre, Mehmet, et al. "Translating C to safer Rust." Proceedings of the ACM on Programming Languages 5.00PSLA (2021): 1-29.
 [2] Zhang, Hanliang, et al. "Ownership guided C to Rust translation." International Conference on Computer Aided Verification. Cham: Springer Nature Switzerland, 2023.
 [3] Eniser, Hasan Ferit, et al. "Towards translating real-world code with LLMs: A study of translating to Rust." arXiv preprint arXiv:2405.11514 (2024).
 [4] Yang, Aidan ZH, et al. "Vert: Verified equivalent rust transpilation with few-shot learning." arXiv preprint arXiv:2404.18852 (2024).

Limitations of existing work



* Correct under tests.

[1] Emre, Mehmet, et al. "Translating C to safer Rust." Proceedings of the ACM on Programming Languages 5.00PSLA (2021): 1-29.
 [2] Zhang, Hanliang, et al. "Ownership guided C to Rust translation." International Conference on Computer Aided Verification. Cham: Springer Nature Switzerland, 2023.
 [3] Eniser, Hasan Ferit, et al. "Towards translating real-world code with LLMs: A study of translating to Rust." arXiv preprint arXiv:2405.11514 (2024).
 [4] Yang, Aidan ZH, et al. "Vert: Verified equivalent rust transpilation with few-shot learning." arXiv preprint arXiv:2404.18852 (2024).











How we conducted our user study

Undergraduates in a computer security course.

Know memory safety and Rust principles.

33 out of 71 students participated.



How we conducted our user study



How we conducted our user study













- > 8 self-contained programs.
- > 300 − 600 LoC.
- Functionality: file/data processing, strings/numbers computation, and parsing.



C to safe Rust is difficult! Translating line-by-line doesn't work

C to safe Rust is difficult! Translating line-by-line doesn't work

Snippet of C benchmark

```
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
  if (condition 1) { *p2 = ' '; }
  if (condition 2) { p1 += 1; }
  p2++;
}
puts(p1);
```

C to safe Rust is difficult! Translating line-by-line doesn't work

Snippet of C benchmark

```
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
  if (condition 1) { *p2 = ' '; }
  if (condition 2) { p1 += 1; }
  p2++;
}
puts(p1);
```





```
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
  if (*p2 == '\t') { *p2 = ' '; }
  if (leading space) { p1 += 1; }
  p2++;
}
puts(p1);
```

Snippet of C benchmark

```
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
  if (*p2 == '\t') { *p2 = ' '; }
  if (leading space) { p1 += 1; }
  p2++;
}
puts(p1);
```

```
unsafe {    Compiler-based translation by Laertes
let mut p1: *mut std::os::raw::c_char;
let mut p2: *mut std::os::raw::c_char; // p1 = ...
p2 = p1;
while *p2 != '\0' as i8 {
    if *p2 == '\t' as i8 { *p2 = ' ' as i8; }
    if leading space { p1 = p1.offset(1); }
    p2 = p2.offset(1);
}
puts(p1); }
```

Snippet of C benchmark

```
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
  if (*p2 == '\t') { *p2 = ' '; }
  if (leading space) { p1 += 1; }
  p2++;
}
puts(p1);
```

```
unsafe {
    Compiler-based translation by Laertes
    let mut p1: *mut std::os::raw::c_char;
    let mut p2: *mut std::os::raw::c_char; // p1 = ...
    p2 = p1;
    while *p2 != '\0' as i8 {
        if *p2 == '\t' as i8 { *p2 = ' ' as i8; }
        if leading space { p1 = p1.offset(1); }
        p2 = p2.offset(1);
    }
    puts(p1); }
```

```
let mut p1: String; // p1 = ...
p1 = p1.replace('\t', " ");
p1 = p1.trim_start().to_string();
println!("{}", p1);
```



```
Snippet of C benchmark
char *p1, *p2; // p1 = ...
p2 = p1;
while (*p2 != '\0') {
    if (*p2 == '\t') { *p2 = ' '; }
    if (leading space) { p1 += 1; }
    p2++;
  }
puts (p1);
```



Difference 1: translation of aliasing pointers



Compiler-based approaches: Line-by-line and pointer-by-pointer



Difference 1: translation of aliasing pointers



Compiler-based approaches: Line-by-line and pointer-by-pointer



Strategy (a): **Elide** aliasing pointers with Rust methods.

```
let mut p1: String; // p1 = ...
p1 = p1.replace('\t', " ");
p1 = p1.trim_start().to_string();
println!("{}", p1);
```

Difference 1: translation of aliasing pointers



Compiler-based approaches: Line-by-line and pointer-by-pointer



Strategy (a): **Elide** aliasing pointers with Rust methods.

```
let mut p1: String; // p1 = ...
p1 = p1.replace('\t', " ");
p1 = p1.trim_start().to_string();
println!("{}", p1);
```

Strategy (b): **Clone** the object to separate r/w access.



Difference 2: translation of C pointers and C API calls

Snippet of C benchmark	Compiler-based translation by Laertes	User translation
	unsafe {	
<pre>char *p1;</pre>	<pre>let p1: *mut std::os::raw::c_char;</pre>	<pre>let p1: String;</pre>
<pre>puts(p1);</pre>	<pre>puts(p1);</pre>	<pre>println!("{}", p1);</pre>
	}	

Difference 2: translation of C pointers and C API calls

Snippet of C benchmark

char *p1; ... puts(p1); Compiler-based translation by Laertes

unsafe {

...

let p1: *mut std::os::raw::c_char;

puts(p1);

User translation

let p1: String;

•••

println!("{}", p1);

 Limited Rust types or unsafe raw pointers.
 unsafe C API calls

10

Difference 2: translation of C pointers and C API calls

Snippet of C benchmark

char *p1; ... puts(p1);



unsafe {

...

let p1: *mut std::os::raw::c_char;

puts(p1);

User translation

let p1: String;

•••

println!("{}", p1);

 Limited Rust types or unsafe raw pointers.
 unsafe C API calls

safe Rust types
safe Rust methods.

Users semantically lift low-level pointers and API calls

Users semantically lift low-level pointers and API calls

User translations:

- \succ Various Rust types.
- \succ No raw pointers.



Users often use zero-cost abstractions



Users semantically lift low-level pointers and API calls

Compiler-based work:

> Call unsafe C APIs

Users:

- ➢ Find safe Rust alternatives.
- \succ Or emulate them.

Direct Mapping (36.7%): Sisdigit Sis_ascii_digit	Map to Expr (39%): fopen(file, "w+") File::options() .read(true)	 Direct Mapping Map to Expression API Emulation Eliminated 3rd party APIs
API Emulation (13.4%): sscanf sscanf_before, sscanf_after	<pre>.create(true) .truncate(true) .open(file)?;</pre>	7.6% 13.4% 36.7%
Eliminated (7.6%): malloc, free, fclose	3rd party APIs (3.3%): geteuid nix::unistd::get euid	39.0%

C features	
Mutable globals	
C unions	



C features	Compiler-based translation	Oser translation
Mutable globals	<mark>unsafe</mark> plain globals	Locals Dynamic references with runtime overhead. Atomic types (lock-free)
C unions	unsafe Rust unions	

C features	Compiler-based translation	C User translation
Mutable globals	<mark>unsafe</mark> plain globals	Locals Dynamic references with runtime overhead. Atomic types (lock-free)
C unions	unsafe Rust unions	Use safe Rust APIs Rust enum

User translations: memory-safe with good performance

User translations: memory-safe with good performance

The execution overhead is mostly within 20%.



Dashed line: the execution time for C programs.

Safe Rust translation: memory safety and good performance is achievable.

Challenging subproblems to address for future automatic translators



LLMs-based tools: decompose the code based on function dependencies.

LLMs-based tools: decompose the code based on function dependencies.

> Inconsistency issue



g2 and g2' are inconsistent!

LLMs-based tools: decompose the code based on function dependencies.

- > Inconsistency issue
- Complicated dependencies



g2 and g2' are inconsistent!



The main of a binary depends on all functions!

LLMs-based tools: decompose the code based on function dependencies.

- > Inconsistency issue
- Complicated dependencies



g2 and g2' are inconsistent!

A better decomposition strategy?



The main of a binary depends on all functions!

For example,

A user translation of fmt

```
1 fn center_stream<R: BufRead>(mut stream: R, _name: &str, config:
         &Config) {
2
       let mut p1 = String::new();
3
       while let Ok(bytes_read) = get_line(stream, &mut p1) {
4
          if bytes_read == 0 { break; }
5
          let len: usize = p1.trim().chars().map(|c| if c == '\t'
        { ' ' } else { c }).map(char::len_utf8).sum();
6
         let padding = (config.goal_length - len) / 2; //
        Calculate padding to center the line
7
8
          for _ in 0..padding { print!("_"); }
           println!("{}", p1.trim());
9
       3 . . .
```

The original C fmt

```
1 static void center_stream(FILE *stream, const char *name)
2 {
3 char *p1, *p2; // aliased pointers
4 size_t len; int w1, w2; wchar_t wc;
5 // ...
6 pair/s (const of the state of th
       6
                                                         while ((p1 = get_line(stream)) != NULL) {
       7
                                                                                  len = 0:
                                                                                    for (p2 = p1; *p2 != '\0'; p2 += w2) {
                                                                                                           if (*p2 == '\t')
         9
       10
                                                                                                                                       *p2 = ''';
       11
                                                                                                            // ... skipped
     12
13
14
15
                                                                                                            if (len == 0 && iswspace(wc)) p1 += w2;
                                                                                                            else len += w1;
                                                                                    3
                                                                                    while (1 < goal_length) {</pre>
       16
17
18
                                                                                                              putchar('_');
                                                                                                              len += 2:
                                                                                    }
 10
19
20
21 }
                                                                                    puts(p1);
                                                       } ...
```

For example,

A user translation of fmt



For example,

A user translation of fmt



For example,

}

puts(p1); } ...

A user translation of fmt 1 fn center_stream<R: BufRead>(mut stream: R, _name: &str, config: &Config) { 2 let mut p1 = String::new(); 3 4 Logical difference! 5 '\t Stdout of Rust: 6 Calculate padding to center the line 7 for _ in 0..padding { print!("_"); } 8 println!("{}", p1.trim()); (empty) Fuzz test 9 CMD: ./fmt -w 10 -c The original C fmt Stdin: \x7a\xc3 X Different behavior! 1 static void center_stream(FILE *stream, const char *name) 2 { 3 4 char *p1, *p2; // aliased pointers size_t len; int w1, w2; wchar_t wc; Stdout of C: 5 while ((p1 = get_line(stream)) != NULL) { 6 len = 0: for (p2 = p1; *p2 != '\0'; p2 += w2) { if (*p2 == '\t') *p2 = '..': 11 // ... skipped 12 13 14 15 if (len == 0 && iswspace(wc)) p1 += w2; else len += w1; while (1 < goal_length) {</pre> 16 17 18 putchar('_'); len += 2:

"Correct" user translations are not fully equivalent to C

"Correct" user translations are not fully equivalent to C



"Correct" user translations are not fully equivalent to C



"Correct" user translations are not fully equivalent to C



"Correct" user translations are not fully equivalent to C

No user translation is fully equivalent.



Specify what behaviors to keep across languages.

We are launching an automatic translation service!

Basic Plan: 5 business days \$100



Premium Plan:

1 business day \$200





Scan to learn more!

Automated Program Translation KISP Lab @NUS

* This slide is for entertainment purposes only.

Breakdown of behavioral differences found by fuzz tests

The best shoco translation





Averaged on translations of all programs

