

Emacs regex compilation

future directions for expressive pattern matching

Danny McClanahan

EmacsConf 2024

Outline

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

Who are you?

Danny McClanahan:

- failed an independent study course in undergrad attempting to make a C compiler because (as my prof *specifically* warned me against!) I got stuck on the parser.
- spent the next several years realizing that was actually what I should have been doing the whole time.

Why are you here?

I have a lot of feelings about text and I'm making it everyone else's problem!

- spent a lot of time learning much of this from emacs-devel 🤔
 - was super confused about... a lot 😓
- thought this would be a matter of just updating the regex engine to use modern techniques 🌀 ⚙️ 🛠️
 - ... turns out the emacs regex engine has features that don't exist in other engines!¹ 🔥 🧠
- then i learned about other larger goals for the regex engine! 🌱
 - ... which happened to overlap a lot with my own research interests 🐱

¹especially non-contiguous input, syntax-aware matching, and multibyte encoding

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

What is a regular expression?

No one's really sure!

When and how does implementation match formal theory?

Formal theory is mostly invoked to post-hoc justify design decisions instead of expanding expressive power!

²This section is adapted from “The Four Eras of Regex” by Prof. Jamie Jennings at NCSU: <https://jamiejennings.com/posts/2021-09-23-dont-look-back-2/>.

How did this happen?

Regex began as an investigation into theories of *formal languages*.

- Many features were added on to implementations to improve the practical user experience.
 - While the people adding these features were often academics, they were still very interested in building practical tools!
- This led to feature development beyond the range of formal theory (awesome!).
 - However, this also means that functionality becomes poorly-specified and implementation-defined.
 - **Formal theory couldn't keep up!**

The 1980s: Moonwalking and Backtracking

- 1983 Michael Jackson demonstrates the *moonwalk* in the music video for "Billie Jean".
- 1986 *Backtracking* is developed to simulate egrep-style regular expressions.

Backtracking Away from Formal Theory

- Backtracking ends up being a fantastic way to implement many more features!
 - ... but this is really where the break from formal theory begins.
- Perl in particular adds a whole host of functionality!
 - ... but this locks people into the specifics of perl as a language to perform many text manipulation tasks.³

³Much bioinformatics code still uses perl!

Formal theory remains largely concerned with incremental improvements to artificial benchmarks, and much less with *expanding* models to cover actual user needs.⁴

⁴This means they're also much slower than they could be if they *listen closely*.

Backtracking Away from Backtracking

By the 1990s, non-backtracking engines are created.

- Extant ones include RE2⁵, hyperscan⁶, and rust regex⁷.
- These make use of the earlier automaton models with linear runtimes for well-specified search tasks.
 - However, they intentionally do not cover anything beyond regular linguistic complexity.

⁵See <https://docs.rs/re2> for the only rust wrapper that will ever be blessed by its maintainer.

⁶See <https://docs.rs/vectorscan-async> for the only rust wrapper that builds it for you.

⁷Of which more will be said later.

So what happens if you need to do more?

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

What are regexps used for?

All variety of text search and parsing tasks!

How do Emacs users use them?

As an auxiliary form of logic, to construct the **user-level grammar for human thought** that Emacs provides: text as input and output.

Why is Text Powerful?

[1: [Text as I/O](#)]

The reason text programming languages are successful is because text is both input (readable) and output (writable).

- This makes text an extremely empowering and accessible framework to *navigate* and *manipulate* program code.

Why is Text Powerful?

[2: [Hidden Dependencies](#)]

If you are unable to modify or deploy your code without employing an opaque external system, then you have a hidden dependency.

- This opaque external system can then exert arbitrary control over your programming output.

Why is Text Powerful?

[3: [Locality](#)]

If you cannot reproduce a system *locally*, it becomes an opaque external system.

- e.g. GUI IDEs, cloud services, Large Language Models⁸

⁸all Microsoft products

Why is Text Powerful?

[4: [What is Text?](#)]

Text is local.

- Emacs is a text editor which implements much of its own logic and user interface via text.
 - This is why we have elisp, a language tightly integrated with text operations from the editor.
- Because text forms UI, *parsing* and *text search* can be employed not just to edit code, but to construct a user interface from text input.
 - This means that language-level mechanisms for text such as the regex engine can be extended into the user interface.

Who Says Text is Empowering?

[1: [You're Not Smart Enough](#)]

Not everyone thinks text is empowering!

- Formal theory thinks nobody should be allowed to parse text without their tools!

Who Says Text is Empowering?

[2: [Don't Parse HTML with Regex](#)]

"Everyone knows" not to parse HTML with regex, because regex (alone!) isn't sufficiently powerful to parse HTML. But:

- 1 Nobody is parsing HTML with a single massive regex!
- 2 Regex + mutable state can achieve arbitrary linguistic complexity!
- 3 Regex search for a specific substring is much faster than parsing everything up front!

Who Says Text is Empowering?

[3: [C Lexer Hack](#)]

Turns out those tools aren't *too* powerful, they're actually not powerful *enough* for practical inputs!

- This is why I got stuck on the parser in that independent study course!

Emacs Says So!

This isn't remotely a concern for Emacs code, which regularly uses regexps to parse HTML and other programming languages! How?

`text properties` write state to the text itself⁹

`syntax parsing` regex engine is aware of this via syntax classes¹⁰

`jit-lock-mode` use smart heuristics to only reparse what's changed¹¹

⁹Not unlike the tape of a turing machine!

¹⁰`\b`, `\<`, etc.: see https://www.gnu.org/software/emacs/manual/html_node/elisp/Regexp-Backslash.html.

¹¹This might just be fontification, as opposed to the work done in `syntax-ppss`.

Formal Theory: Right for the Wrong Reasons

There *are* actually reasons to avoid using regexps to parse text!

- Regexps may have extremely non-obvious dependencies on parse context.
 - A non-greedy match may be correct when invoked in a restricted context, but may become subtly incorrect when used more generally.¹²
- While text properties and buffer-local variables can retain the state necessary to parse non-regular languages, coordinating that state can be error-prone.
 - Since **there are no existing formalisms to link regex with external state**¹³, it can become extremely difficult to reproduce the precise internal state which generates a logic bug in an elisp mode.

¹²For example, `(_<.*?)`: could match a symbol before a `:` (like `a:` in JavaScript), but could unintentionally match string properties like `"a:b"`: as `a:` too!

¹³Composing automata with other parse state is one of the subjects of my research.

In fact, `tree-sitter` (since Emacs 29) was created to solve this problem *for well-specified language definitions*.

- It is a highly constraining formal tool!
- And it means you now depend on:
 - The `tree-sitter` grammar for your language.¹⁴
 - The `tree-sitter` library.¹⁵
- So I don't like it!
 - But for the specific task of parsing a programming language, it happens to solve a lot of other problems at once.

¹⁴obnoxious to read and write

¹⁵does not have universal uptake within distros

So Why Use Regex?

So why are we talking about regex here? Mainly:

- Parsing programming languages is a very small subset of all text search/matching tasks!
- **Regex can be directly manipulated by the user!**

For the interactive experiences that Emacs excels at, regex provides a powerful language *for both input and output*:

- It can be synthesized hygienically from elisp code via `rx`, either statically at load time or dynamically at run time!
- It can be received or transformed from user input to specify powerful queries over complex data!¹⁶

¹⁶See `helm-rg` and `telepathygrams` at end.

...but this might require going beyond "regex" alone!

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?**
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

What is the Emacs regex engine?

It's a backtracking engine over multibyte codepoints, defined in `src/regex-emacs.c`.

How is it invoked?

In two ways:

- over a single contiguous string input,
- over the two halves of the gap buffer.

¹⁷This section is an unfortunately brief walkthrough through the current regex logic.

The compiled pattern is stored as an `re_pattern_buffer` struct from `src/regex-emacs.h`.

- In particular, unsigned `char *buffer` holds the instructions!
- Case folding uses the char table in `Lisp_Object translate`.

The matching loop in `re_match_2_internal()` in `src/regex-emacs.c` goes vaguely as follows:

- 1 extract current and next char
 - perform multibyte varint decoding to iterate bytes
 - translate input characters via the `translate` case-folding char-table
- 2 read instruction from instruction pointer
- 3 big switch statement for the next instruction
 - if instruction uses syntax, read the syntax class¹⁸ for the current character from the current syntax table
- 4 increment the instruction pointer¹⁹
- 5 if we've concluded a capture, write the end position to the C-level array `re_nsub`

¹⁸https://www.gnu.org/software/emacs/manual/html_node/elisp/Syntax-Class-Table.html

¹⁹unless instruction was a jump

Non-Contiguous Matching

Non-contiguous matching over the two halves of the gap buffer is supported by checking at each point whether we have progressed to the end of the first half, and then switching over to the second half.

- This allows the same code to be used for single-string search, as it simply avoids checking a NULL second pointer and only checks if we've reached the end of the first input.

- It turns out this actually isn't terribly relevant to the regex engine!
 - Or at least, it doesn't really differ from "standard"²⁰ Unicode regex matching.
- Emacs reads in data from whatever encoding into multibyte,²¹ and the regex engine only acts upon this normalized encoding.

²⁰There is no (real) standard (yet):

<https://jamiejennings.com/posts/2021-09-07-dont-look-back-1/>.

²¹See https://www.gnu.org/software/emacs/manual/html_node/elisp/Text-Representations.html.

Is that all?

How much time²² do you have?

²²and space

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

How could we do regex better in Emacs?

- introspection
- optimization

How could Emacs do regex better than anywhere else?

- explicit control over linguistic complexity
- libraries of composable patterns

²³This section will describe several potential paths we might investigate, paraphrasing discussion from `emacs-devel` .

Separately-Compiled Regexps

Precompile regexps to enable more powerful compilation techniques.

Problem Emacs currently uses a fixed-size global compile cache.

Solution Create native elisp objects for regexps and match data.²⁴

- Use `(make-regexp "...")` to explicitly compile a pattern string.
- All supported²⁵ methods in `search.c` can accept either a compiled regexp or bare string.

²⁴I have demonstrated this in a test branch:

<https://github.com/cosmicexplorer/emacs/tree/lisp-level-regex>.

²⁵Literal search methods do not use the regexp engine.

Separately-Compiled Regexps

Precompile regexps to enable more powerful compilation techniques.

Results Artificial benchmarks²⁶ show an improvement!²⁷

- ... but I haven't been able to produce an apples-to-apples comparison yet.
- Syntax highlighting would be the most appropriate, but caching these compiles currently breaks fontification.²⁸

²⁶Using `test/manual/perf.el`.

²⁷Using the native match data object produces no improvement over consing a list.

²⁸Current guess is that it relies on buffer-local state not available when precompiled.

Match Over Bytes, not Chars

Compile patterns to byte-level automata, then iterate over bytes.

Problem Char-by-char varint decoding of multibyte²⁹ is comparatively slow.

- This is the reason go's "RE2" is much much slower than the C++ RE2 library.³⁰

Solution We can do this work at compile time instead!³¹

- Generates a larger automaton in order to be able to think in terms of byte ranges.³²

²⁹No worse than UTF-8 in general, but we may be able to pipeline decoding somehow.

³⁰Source: a very dear friend.

³¹This is already what we do for e.g. char-folding.

³²This is a necessary prerequisite for SIMD instructions.

Explicit Control over Linguistic Complexity

Break apart the monolithic regex pattern interface into subroutines for specific inputs.

- Problem[1]** There's no way to validate that a given pattern isn't more complex than expected.
- This requires careful escaping to avoid accidentally triggering regex behavior.
- Problem[2]** There's no way to ensure Emacs uses faster algorithms³³ for less complex patterns.³⁴
- This results in difficult-to-understand performance characteristics.
- Problem[3]** There's no way to specify different search semantics.³⁵
- Instead, we have a single type of input and a single type of output.

³³See <https://github.com/BurntSushi/rebar> for a fantastic discussion of techniques for regex performance.

³⁴We perform a heuristic check for literal patterns, but only in some code paths.

³⁵Such as allowing false positives, or matching against a set of patterns.

Explicit Control over Linguistic Complexity

Break apart the monolithic regex pattern interface into subroutines for specific inputs.

- Solution[1]** Single or multiple literals³⁶ can employ specialized SIMD algorithms to avoid reading every single byte one by one.³⁷
- Solution[2]** Non-capturing patterns or patterns without backrefs³⁸ can use faster automata.
- Solution[3]** Collecting a sequence of matches for the same pattern can be done all at once.³⁹
- Solution[4]** Matching against a set of patterns can be done more efficiently and ergonomically than combining with `\|`.

³⁶Multiple literals is especially helpful for matching a set of keywords.

³⁷This is used as a "prefilter" optimization in modern engines like RE2 to avoid reading each byte one-by-one.

³⁸These have recently been formalized:

<https://jamiejennings.com/posts/2023-10-01-dont-look-back-3/>.

³⁹Overlapping matches can be supported for ambiguous cases, instead of choosing longest or shortest only.

Lisp Regexp Library

Expose a lisp-level library for regexp matching.

Problem The compiled form of the regexp in `re_pattern_buffer` can be *executed*, but not really *introspected*.

- No form of "IR": this also contributes to the difficulty of composing patterns together.

Solution⁴⁰ We have `libgccjit` now: **why not implement the regexp engine itself in lisp?**⁴¹

Results (*postulated*):

- Integration into `pcase` could achieve a form of type safety along with interleaving lisp-level matching logic.
- Biggest issue for optimization: lisp code (or native modules) can't access or operate on the separate halves of the gap buffer.

⁴⁰Proposed by Pip Cet on `emacs-devel`.

⁴¹Alternatively, translate the regexp into lisp which we can then JIT.

Question

- 1 Who are you? Why are you here?
- 2 What is a regular expression? When and how does implementation match formal theory?
- 3 What are regexps used for? How do Emacs users use them?
- 4 What is the Emacs regex engine? How is it invoked?
- 5 How could we do regex better in Emacs? How could Emacs do regex better than anywhere else?
- 6 Do you have any concrete examples? Or are you just posturing?

Do you have any concrete examples?

Yes!

Or are you just posturing?

My posture is terrible!⁴²

⁴²Other translations have also been suggested by modern scholars, including *formidable*, as well as *awe-inspiring*.

helm-rg⁴³

A code search tool similar to M-x grep, using ripgrep⁴⁴.

- Generates regexps from input:
 - "a b" => "a.*b|b.*a".
 - Translates from PCRE to elisp regexps to highlight matches in the helm buffer.
- M-b enters "bounce mode", where matched lines can be edited directly.

Example (pattern generation)

```
return sources

def compute_extra_classpath(self, extra_compile_time_classpath_elements):
    """Compute any extra compile-time-only classpath elements.

    TODO(benju): Model compile-time vs. runtime classpaths more explicitly.
    """
    def extra_compile_classpath_iter():
        for conf in self.configs:
            for jar in extra_compile_time_classpath_elements:
                yield (conf, jar)

    return list(extra_compile_classpath_iter())

@memoized_method
def _plugin_targets(self, compiler):
    """Returns a map from plugin name to the targets that build that plugin."""
    if compiler == 'javac':
        plugin_cls = JavaPlugin

-- L105/L075 97% C15 jvm_compile.py (Python Gitter II hi-p Undo-Tree Sert Rbow RRow 80+ Helm
C-; Visit result buffer and highlight matches (keeping session)
rg argv: ( /home/cosmicexplorer/.cargo/bin/rg --smart-case --color=ansi --color=
file "extra_classpath_compile.*classpath.*classpath.*extra.*compile[classpa
src/python/pants/notes/easter_egg
3402* Put all extra classpath elements (e.g., plugins) at the end (scale compile)
src/python/pants/backend/jvm/tasks/jvm_compile/jvm_compile.py
339: def extra_compile_time_classpath_elements(self):
321:     """Extra classpath elements common to all compiler invocations.
432:     self.extra_compile_time_classpath_elements(),
517:     extra_compile_time_classpath_elements();
540:     extra_compile_time_classpath = self._compute_extra_classpath(
542:     extra_compile_time_classpath_elements)
548:     extra_compile_time_classpath.
548: def _create_compile_job(self, classpath_products, compile_contexts, extra_compile_time_classpath,
513:     extra_compile_time_classpath,
305:     """Compute any extra compile-time-only classpath elements.
305* """Compute any extra compile-time-only classpath elements.
309: def extra_compile_classpath_iter():
309:     for jar in extra_compile_time_classpath_elements:
309:         return list(targets.generated_classpath_iter())
helm-rg 1/1 (C) Candidate(s) C-h #Help TAB/Act RET/1/2/F-m:RkAct C-!;on.suspend
rg pattern: extra_compile_classpath
```

⁴³<https://github.com/cosmicexplorer/helm-rg>

⁴⁴<https://blog.burntsushi.net/ripgrep/>

telepathygrams⁴⁸

A (WIP) code search tool that precompiles a database to execute NFAs against.

I want to beat ripgrep by "cheating"⁴⁵ with a precompiled index.⁴⁶

- n-gram indices have been done,⁴⁷ but I don't want to just find where to *start*—I want to execute the *entire search* against the index!
- This requires virtualizing NFA state so that it may be distributed:
 - *across time* in parallel / across machines,
 - & *space* in terms of *offsets* vs directly against the input data.
- This may fail, but it will be fun!

⁴⁵There are other ways to cheat here too, like precompiling known queries.

⁴⁶Inspired by etags, but with a more complex index for more general queries.

⁴⁷e.g. Kythe: <https://kythe.io/docs/kythe-overview.html>

⁴⁸<https://github.com/cosmicexplorer/telepathygrams>

Paul Wankadia taught me everything.

⁴⁹<https://github.com/google/re2/issues/502>

(point-max)

 call me beep me if you wanna reach me⁵⁰

text-mode

`fed` @hipsterelectron@circumstances.run^a

`IRC` @cosmicexplorer in #emacsconf^b on irc.libera.chat

`email` dmc2@hypnicjerk.ai^c

^aalso (sporadically) twitter & bluesky

^balso (sporadically) elsewhere

^calso dmc2@ (still deciding which poetic license to use)

prog-mode

`codeberg` @cosmicexplorer

`github` @cosmicexplorer

⁵⁰or hire me