

# Interview Notes

Meshal Binnasban

October 21, 2025

# 1 Questions

## **Q1: Why not just use existing regex engines?**

Existing engines are efficient in practice but still face size-related limitations. DFA-based engines such as rust attempt to avoid constructing the full automaton by starting from the initial state and building DFA states only as needed for the given input string. However, for patterns such as  $r^n$ , the engine still needs to produce  $n$  copies of  $r$ 's automaton. Although rust uses size-checking and pattern heuristics to decide which regexes are “safe,” this behaviour is not formally guaranteed, and some cases still slip through, approaching memory limits and running slowly.

Python, on the other hand, uses an NFA and does not generate a DFA. This design supports additional features like backreferences and lookarounds but sacrifices predictability. While an NFA has linear size ( $n + 1$  states), it can still require exploring many or even all states. Python uses depth-first search, which is fast when a match is found quickly, but if not, the engine might need to explore a large number of paths—or even the entire tree—before reaching a decision. If the input is incorrect, the search may still traverse the full tree, leading to catastrophic backtracking.

## **Q2: What happens when we add backreferences?**

Adding backreferences changes the matching problem: equality of unbounded substrings cannot be expressed by any finite automaton, so subset construction does not apply as well. the general matching problem is NP-complete, and practical engines therefore typically rely on backtracking search

## **Q3: What about derivative-based matchers?**

Derivative-based matchers are elegant and easy to reason with, and they have been proven correct for both acceptance and lexing. However, the Achilles' heel of the derivative approach is expression growth. Each derivative reconstructs new parts of the expression, particularly for SEQ and STAR, which can expand rapidly. Even with aggressive simplifications, the size can still grow, since equivalent subexpressions may appear at different levels of the expression tree and cannot be merged.

## **Q4: How does the marked approach help?**

The marked approach avoids expression reconstruction altogether. It keeps the structure of the regex fixed and moves marks through it as matching progresses. These marks represent how far the matcher has advanced inside each subexpression. When a character is read, the marks shift through the regex to indicate progress. This approach eliminates structural duplication and

makes it easier to reason about correctness and control expression growth. Fischer’s original algorithm focused only on acceptance, while my work extends this framework to also construct values, enabling further reasoning about POSIX disambiguation.

**Q5: What exactly does POSIX mean in your work?**

In this work, POSIX follows the definition of Tan and Urban (2023): among all possible matches, the POSIX value is the one that is leftmost and, among those, longest. This selection is made by comparing how early and how completely each possible match consumes the input.

**Q6: Why is value extraction important?**

Value extraction records how the regex matched, not just whether it matched. It is useful in several areas:

- Tokenization: for instance, splitting a source file into identifiers, numbers, and punctuation before parsing.
- Syntax highlighting: regexes are often used to detect keywords, comments, and string literals; reconstructing the match (e.g., a value like `Left(Right(Empty))`) tells the highlighter which span of text to colour.
- Lexing: in compilers, regexes define each token class, and the constructed values help establish token boundaries and precedence.

Accurate and reproducible value extraction ensures that tools relying on regex matching behave deterministically and consistently.

**Q7: What problem requires research?**

Derivative-based matchers face expression size explosion, and existing marked algorithms only handle acceptance. Using marks helps avoid this by keeping the regex structure fixed and instead moving progress markers through it. This prevents the exponential growth of states seen in automata and the structural duplication of derivatives. Marks also provide a clear foundation for value construction, where each mark can carry additional information such as bit annotations or the remaining suffix of the string.

The next step in this line of research is to define an ordering of marks that captures the POSIX leftmost-longest preference directly within the algorithm’s shifting process. This will allow the matcher to produce the POSIX-correct value by construction, rather than by comparing alternatives after matching.

**Q8: What is the difference between your work and Fischer’s?**

Up to this stage, we have developed an algorithm that extends Fischer’s

marked approach so that marks carry strings, beginning with the full input string. As matching progresses, each mark represents how far its branch has consumed the input. This enables both acceptance and value construction. The main difference is that Fischer’s version only checked whether a match existed, while my version records and reconstructs the structure of that match—preparing the ground for later disambiguation according to POSIX rules.

#### (Summary of contribution)

This approach aims to avoid the limitations of state explosion, backtracking, and derivative size growth, while providing not only matching but also value extraction with POSIX disambiguation.

## 2 Notes

### 2.1 proofread notes

String matching is usually done by constructing an NFA from a regular expression, then compiling it into a DFA, and finally matching character by character. This process can be efficient for small patterns but may lead to a state explosion in the DFA, reaching up to  $2^n$  states for expressions like  $r^n$ . Engines that rely on these constructions use various strategies to manage this growth, such as building DFA states lazily or exploring nondeterministic paths, each with its own trade-offs.

#### 2.1.1 Flow

The classical path is regex → NFA → DFA. We construct an NFA from a regex. Then we apply subset construction to obtain a DFA. This “regex → NFA” stage is a helpful stepping stone before determinization.

#### 2.1.2 How string matching is done (DFA route)

We first produce an NFA. We then build a DFA. Finally, we match character by character on the DFA. In a DFA, the running time is naturally linear in the input length: each character triggers a single transition, and we can tell if we end in an accepting state. The main limitation is automaton size.

### 2.1.3 Rust's regex engine (claim and scope)

Rust's engine does not support lookarounds or backreferences. In return, it offers a linear-time guarantee for the regex and the text, but only for approved patterns. With patterns like  $R^n$ , the situation deteriorates: translating  $R^n$  to a DFA effectively chains  $n$  copies of the DFA for  $R$ . The guarantee is tied to a size heuristic. It holds only within the approved set. Some hard cases slip the heuristic, can approach memory limits, and run slowly.

### 2.1.4 Automata size and lazy construction

When generating a DFA, the number of states can reach  $2^n$  in the worst case, with  $n$  proportional to the regex size. Rust does not build a complete DFA up front. It constructs states lazily and on demand, often yielding a small DFA on typical inputs. For adversarial inputs, one may still end up generating essentially all  $2^n$  states.

### 2.1.5 Python/PCRE engines (NFA with backtracking)

Python does not generate a DFA. It executes an NFA with backtracking. A standard construction yields an NFA with size linear in the regex. For  $R^n$ , this means  $n$  sequential copies of the fragment for  $R$ . Because NFAs are nondeterministic, matching may need to explore many paths. Two strategies are common:

- Depth-first search (DFS): follows one path, often fast when acceptance is nearby, but susceptible to catastrophic backtracking.
- Breadth-first search (BFS): explores level by level; in the worst case memory grows exponentially.

A practical reason for this design is feature support: Python/PCRE include backreferences, which preclude a clean DFA route.

### 2.1.6 Backreferences and complexity

Backreferences push matching beyond regular languages. The general matching problem with backreferences is NP-complete. Backreferences enforce equality of arbitrarily long substrings, so there is no general subset construction to a finite DFA.

### 2.1.7 Derivatives

Derivative-based matchers do not support backreferences. They are typically fast in common cases. They are elegant and proof-friendly, with correctness proofs for acceptance and lexing. The Achilles' heel is potential expression growth, especially through SEQ (product). We apply simplification to curb this growth. There are no worst-case guarantees. Duplicated subexpressions can appear at different levels or positions, where local simplifications cannot merge them.

### 2.1.8 Example note (backtracking)

To illustrate catastrophic backtracking in Python-style engines, ( $(a^*) + b$  on an all a input).

## 2.2 raw notes

1 – *Regex – > NFA – > DFA*. constructing from REGEX to NFA is better stepping stone to then convert into DFA. how is string matching done? String matching is done by producing NFA then DFA then match character by character. in DFA, some claim matching is linear to input? at first, it might make sense since you have the string and you always know where to transition in a DFA, to know if you are in an accepting state. rust claims "they dont have lookaround and backreferencing in exchange to linear time with respect to the size of the regex and the search text. however, if you feed it a specific regex such as one with power to n ( $r^n$ ) it doesnt. if we translate regex of a power of n to DFA, it will copy the r DFA into n copies and put them together. rust then says the it is linear for approved regex, limited by size which they determine using an algorithm. they give this strong property but only in the case they deem ok. ofcourse some cases slips thier algorithm, it might not exhaust the memory as some regex do, but just about exhausting it and it is slow. - note : when generating DFA, size can grow  $2^n$ . rust dont generate complete automata (because a small regex could exhaust memory) they generate DFA essentially dynamically (they generate the starting state and depending on the input string, usually relativly small DFA) but of given the incorrect string , you still have to generate all the  $2^n$  states. python why perform bad on example of  $a^{**}b$ ? in python, they dont generate DFA, instead they generate NFA. there is an NFA property

which is that the number of states is  $n + 1$  if you have a regex with size  $n$  in the worst case, unfortunately you will have to generate  $n$  copies for  $r^n$ . if the matching is done in NFA, you might have to explore all states to find accepting states. because, in nfa, each state you don't know where exactly to go unlike DFA. to implement this (meaning matching regex using nfa), you can use 1- depth first search or 2- breadth first search. 1- in depth, explore one path and count on finding accepting state very quickly, which is fast most of the time. 2- in breadth search, to explore every level of the tree at the same time, sometimes you may need to explore the full tree, and the size grow exponentially for the memory. Downfall of python, they use this to include back-referencing, so that is why they use nfa, downfall is if the string don't match, you may need to explore the full tree, tricking the engine with depth first causes this. note- if backreference is added , then the matching problem could become np complete problem ( or equivalent). most efficient algorithms for solving np problems are exponential algorithm (maybe an example such as backtracking heuristic algorithm), backreference changes the problem and there is no possibility of subset construction (explain why here). how about derivative: no backreferencing. fast in usual (such as examples given in lecture ) but also has cases where it can explode in size. achilles heel of derivative that it can grow quickly and make it slower. to try and prevent this, simplification is done. unfortunately no guarantees that simplification works all the time, copies to be simplified might not be possible to find, if  $r$  and its copies are in different levels and positions in the tree and next to eachother, then simplification cannot be applied. as elegant and beautifully designed derivatives are, and easy to reason about and implement (and proofs exist of it and its lexer as well which is something some other matcher implementations lacks), it still has this problem.