

supported by:



A dependently-typed regex parser.

Ohad Kammar and Katarzyna Marek

University of Edinburgh

```
r "(([01][0-9])|([2][0-4])):([0-5][0-9])"
```

	Validation : Bool	Parsing : Maybe (Int, Int)
"12:03"	True	Just (12, 3)
"12:63"	False	Nothing

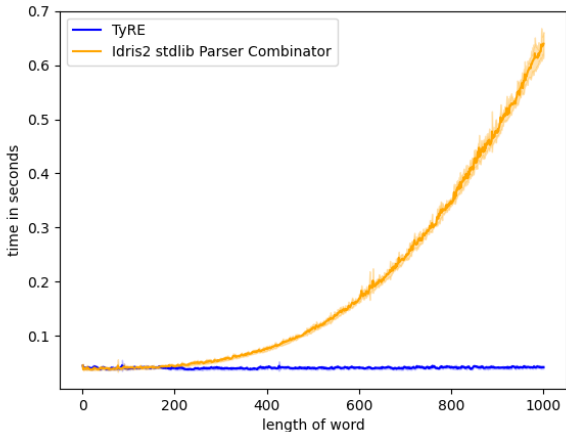
G. Radanne. [PEPM'19]

Typed parsing and unparsing for untyped regular expression engines.

Idris-TyRE vs. Radanne's tyre:

- tyre layer - *digit* <*> *any* : *TyRE* (*Int*, *Char*)
- unparsing - *Conv snd (r "[AB][0-9]")*
- custom parser (there is no regex parser in Idris)
- type safety throughout all layers

WHY REGEXES?



regex: $((a*c)|a)*b$; string: $a^{n-1}b$

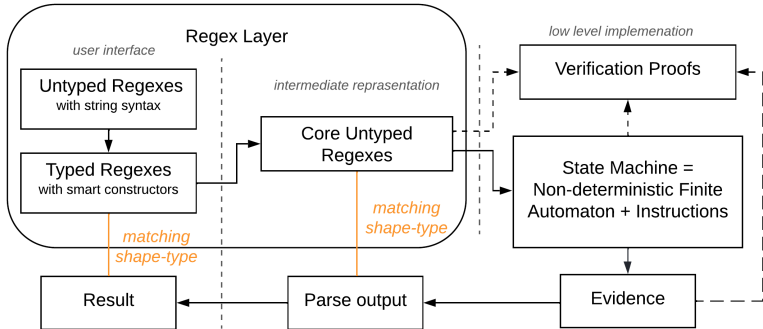
SHAPE OF A REGEX

r "[AB][0-9]" : TyRE (Char, Char)
shape

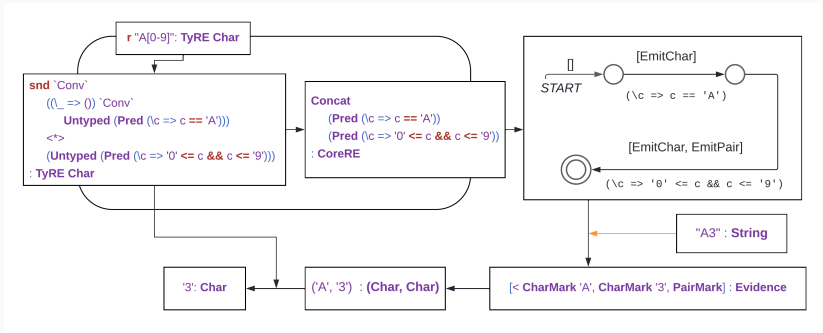
Literal	RE	Shape
<i>a</i>	<i>Exactly a</i>	<i>Unit</i>
<i>[abc]</i>	<i>OneOf ['a', 'b', 'c']</i>	<i>Char</i>
<i>[a-c]</i>	<i>'a' `To` 'c'</i>	<i>Char</i>
<i>.</i>	<i>Any</i>	<i>Char</i>
<i>AB</i>	<i>A `Concat` B</i>	<i>(Shape A, Shape B)</i>
<i>A B</i>	<i>A `Alt` B</i>	<i>Either (Shape A) (Shape B)</i>
<i>A?</i>	<i>Maybe A</i>	<i>Maybe (Shape A)</i>
<i>A*</i>	<i>Rep0 A</i>	<i>List (Shape A)</i>
<i>A+</i>	<i>Rep1 A</i>	<i>List (Shape A)</i>
<i>`A`</i>	<i>Group A</i>	<i>String</i>

Demo

ARCHITECTURE

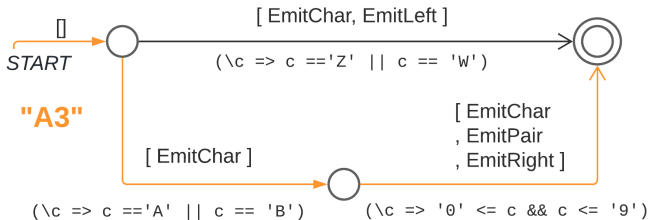


EXAMPLE



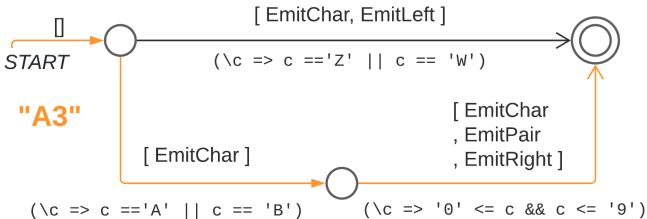
VERIFICATION PROOFS I - CREATING AN ACCEPTING PATH

Given a state machine SM, that accepts a string S,
then there exists an accepting path through SM,
s.t. the evidence collected by following this path is equal to evidence
collected when running the SM on S.



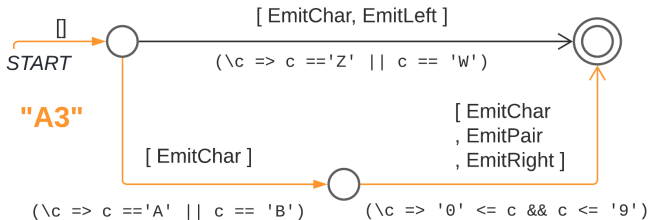
VERIFICATION PROOFS II - EVIDENCE SHAPE

Given a state machine SM, that was constructed for a regex RE
and an accepting path A through SM,
then the evidence collected by following A through SM encodes the
shape of RE.



$r \text{ "[ZW]|([AB][0-9])"}$
: TyRE (Either Char (Char, Char))

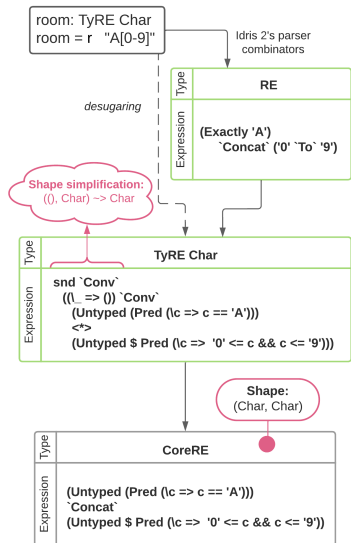
r "[ZW]|([AB][0-9])"
: TyRE (Either Char (Char, Char))



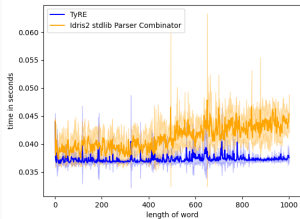
The proofs are **erased at runtime**.

STRING SYNTAX

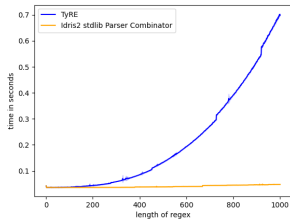
1. We parse strings into *RE* during compile-time.
 - We use a copy of Idris 2's parser combinators library.
2. More type-level calculation to get the *TyRE*.
 - One to one translation.
 - Shape simplification.



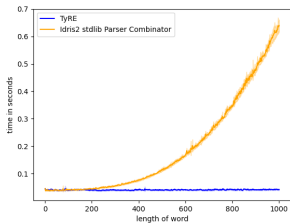
PERFORMANCE



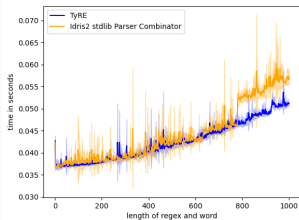
regex: a^* ; string: a^n



regex: $a(|a)^{n-1}$; string: a



regex: $((a^*c)|a)^*b$; string: $a^{n-1}b$



regex: a^n ; string: a^n

LIMITATIONS

- Library supports only regular expressions (in formal language theory meaning).
- Minimal performance engineering
 - no support for SM pre-compilation,
 - no automata minimisation.
- We prove safety, not functional correctness.
- Missing features:
 - regex string literals syntax, e.g. "a{0, 3}" is not supported,
 - common regex usages, e.g. find all and replace.
- Compile time performance:
 - parsing string literals,
 - running TyRE at compile-time.

- Integrate TyRE with Idris 2's stdlib parser combinators library.
 - This would allow bootstrapping and using TyRE to parse regex string literals.
 - Challenge : compile-time parsing.
- Adding unparsing.
 - Possible use case:
Collie: command-line interface generator
- Extending to μ -regular expressions (context-free expressions) [Krishnaswami and Yallop, 2019].



[Allais and Kammar]

Questions?

REFERENCES



Krishnaswami, N. R. and Yallop, J. (2019).

A typed, algebraic approach to parsing.

In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019*, page 379–393, New York, NY, USA. Association for Computing Machinery.



Radanne, G. (2019).

Typed parsing and unparsing for untyped regular expression engines.

In *PEPM 2019 - Proceedings of the 2019 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation, Co-located with POPL 2019*.