

# Search-Based Regular Expression Inference on a GPU

Mojtaba Valizadeh

Martin Berger

<https://martinfriedrichberger.net/>

21 June 2023

# Problem

Can we accelerate program synthesis with GPUs?

## What makes program GPU-friendly?

- ▶ Predictable data movement
- ▶ Minimise data-dependent branching
- ▶ Maximise parallelism
- ▶ Minimise synchronisation between different threads/warps/cores etc as much as possible

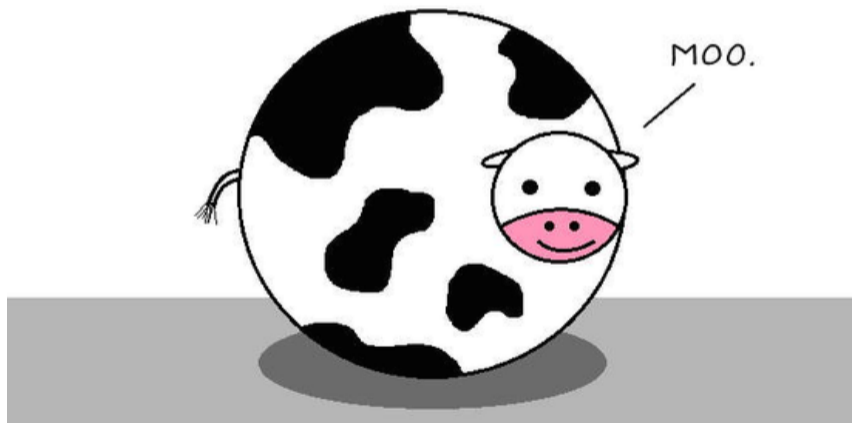
Can we do this for program synthesis?

## What makes program GPU-friendly?

- ▶ Predictable data movement
- ▶ Minimise data-dependent branching
- ▶ Maximise parallelism
- ▶ Minimise synchronisation between different threads/warps/cores etc as much as possible

Can we do this for program synthesis?

What is a tractable synthesis problem to bootstrap GPU synthesis?



Regular expression inference

REI = programming by example with regular expressions

# REI = programming by example with regular expressions

## Input:

- ▶ **Examples:** two finite sets of strings,  $P$  and  $N$
- ▶ **Cost function:**  $\text{cost}(\cdot)$  for REs

## Output: regular expression $r$ that is:

- ▶ **Precise:**  $r$  accepts all strings in  $P$  and rejects all strings in  $N$
- ▶ **Minimal:** no regular expression with a cost less than  $\text{cost}(r)$  is precise



## REI = programming by example with regular expression

Old & well-known problem (1967 Gold's "Language identification in the limit")

Currently unsolved for deep-learning based

REs are "embarrassingly sequential"

Asymptotic complexity well understood

- ▶ Gold (1978): is NP-hard for DFAs
- ▶ Pitt and Warmuth (1993): is NP-hard for DFAs, NFAs and REs, even to approximate minimum

# Regular expressions

The **regular expressions** over  $\Sigma$  are given by the following grammar:

$$r ::= \emptyset \mid \epsilon \mid a \mid r \cdot r \mid r + r \mid r^*$$

We write  $\text{Lang}(r)$  for the language of  $r$

Usual abbreviations, e.g.,  $rr'$  for concatenation  $r \cdot r'$ ,  $r?$ ,  $\Sigma^*$ . Also implemented: intersection, negation, complement

## Search order, cost homomorphism

A **cost function** is a map  $\text{cost}(\cdot) : \text{RE}(\Sigma) \rightarrow \text{Nat}$ . It is a **cost homomorphism** if there are integer constants,  $c_1, \dots, c_5 > 0$ , such that:

- ▶  $\text{cost}(\emptyset) = \text{cost}(\epsilon) = \text{cost}(a) = c_1$  for all  $a \in \Sigma$
- ▶  $\text{cost}(r?) = \text{cost}(r) + c_2$
- ▶  $\text{cost}(r^*) = \text{cost}(r) + c_3$
- ▶  $\text{cost}(r \cdot r') = \text{cost}(r) + \text{cost}(r') + c_4$
- ▶  $\text{cost}(r + r') = \text{cost}(r) + \text{cost}(r') + c_5$

We call each  $c_i$  the **cost** of the corresponding regular constructor

## Solution: trivial algorithm

Enumerate by increasing cost and check each if meets constraints

## Solution: trivial algorithm

Enumerate by increasing cost and check each if meets constraints

**Research question:** How to make GPU friendly?

## Core problem: how to represent REs during search?

Using REs search space is wasteful for several reasons:

- ▶ **Not space-efficient:** Each regular expression is a tree (needs pointers)
- ▶ **Slow contains-check:** To check candidate, needs tree-walk
- ▶ **Redundant:** Each regular language is denoted by infinitely many regular expressions. For example,  $00+1$  and  $1+00$  denote the same language
- ▶ **Non-local:** Pointers can point to anywhere in memory, unpredictable

## Representation of RE by languages

Recall language  $L$  is subset of  $\Sigma^*$ . Isomorphic representation as **characteristic sequence**

$$\mathbf{1}_L : \Sigma^* \rightarrow \mathbb{B}$$
$$\sigma \mapsto \begin{cases} 1 & \sigma \in L \\ 0 & \text{else} \end{cases}$$

Note:  $\mathbb{B}$  is a semiring, and boolean algebra

Many REs have the same characteristic sequence e.g.  $r + r$  and  $r$ . Characteristic sequence prune away this redundancy (REs up to observational congruence)

Problem: the domain  $\Sigma^*$  of  $\mathbf{1}_L$  is infinite

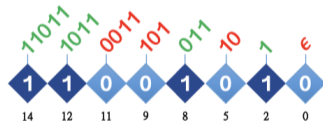
# Representation of languages

In REI, we only care about strings in  $P \cup N$ , we could use

$$\mathbf{1}_L : P \cup N \rightarrow \mathbb{B}$$

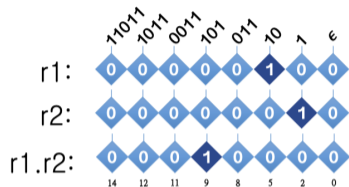
Positive = {1, 011, 1011, 11011}

Negative = { $\epsilon$ , 10, 101, 0011}





Note:  $\mathbf{1}_L : P \cup N \rightarrow \mathbb{B}$  is not enough, e.g. with



$P = 1, 001, 1001, 11011$

$N = \epsilon, 10, 101, 0011$

$1011 = \text{"".}1011$  checked  
 $= 1.011$  checked  
 $= 10.11$  **!?**  
 $= 101.1$   
 $= 1011.\text{""}$

## Infix-closure

Our representation of regular expressions are characteristic sequences

$$cs : ic(P \cup N) \rightarrow \mathbb{B}$$

$w$  is an **infix** (aka substring) of string  $\sigma$  if  $\sigma = w_1 \cdot w \cdot w_2$ . Example: “abc” has infixes like

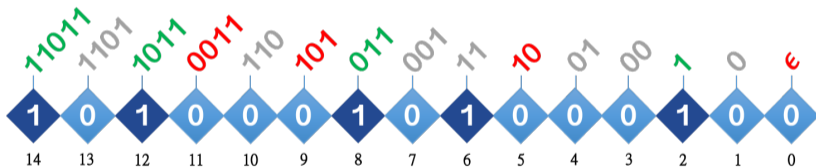
*abc ab bc a b c  $\epsilon$*

$ic(S)$ , the **infix-closure** of a set  $S$ , is the smallest infix-closed superset of  $S$

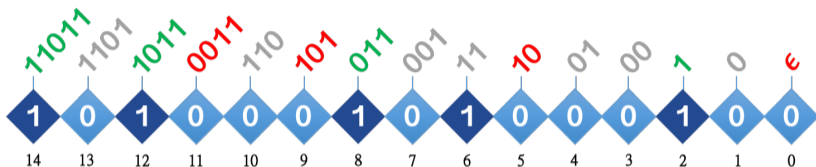
## Infix-closure: example

$P = \{1, 011, 1011, 11011\}$  and  $N = \{\epsilon, 10, 101, 0011\}$ . Then  $\text{ic}(P \cup N)$  is

11011, 1101, 110, 11, 1011, 101, 10, 1, 011, 01, 0011, 001, 00, 0,  $\epsilon$



## Infix-closure characteristic sequences summary



- ▶ Characteristic sequence turns language into a bitvector in memory
- ▶ bitvector = integer, bitvectors = integer matrix!
- ▶  $P$  and  $N$  are fixed, so don't change during an REI run!
- ▶  $ic(P \cup N)$  is ordered in memory
- ▶ Going through  $ic(P \cup N)$  is a linear scan (= predictable)

# Principled program synthesis: using formal power series over semirings

Our characteristic sequences  $cs : ic(P \cup N) \rightarrow \mathbb{B}$  form a **\*-semiring**, abstracting:

- ▶ addition
- ▶ multiplication
- ▶ iteration

corresponding exactly to the operations on regular expressions

**Slogan:** \*-semiring as the API for regular expression synthesis

---

### Algorithm 1 Main function of synthesis algorithm

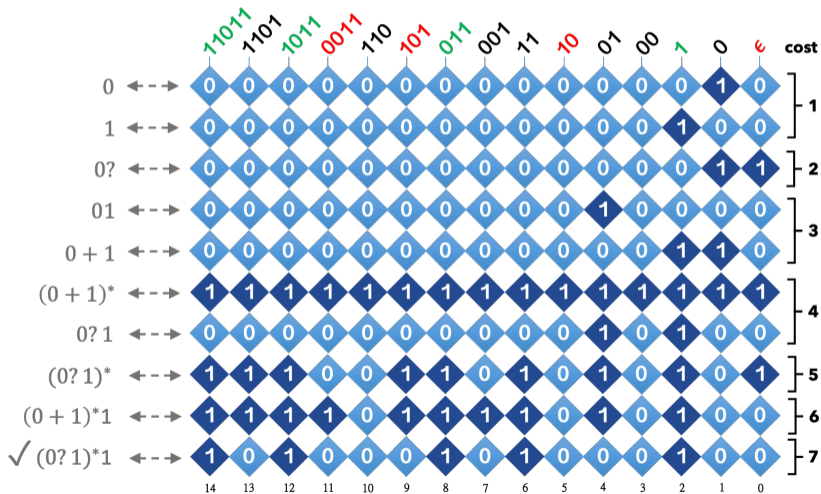
---

1: **Input** Positive and negative examples  $(\mathcal{P}, \mathcal{N})$ ,  $cost$ ,  $maxCost$   
2: **Output** A minimal RE w.r.t. cost and  $maxCost$  and consistent with  $(\mathcal{P}, \mathcal{N})$ , otherwise "not\_found"  
3:  
4: **if**  $\mathcal{P} == \{\}$  **then return**  $\emptyset$   
5: **if**  $\mathcal{P} == \{''\}$  **then return**  $\epsilon$   
6:  $languageCache =$  [list of CSs of alphabet] ▷  $languageCache$  is global variable  
7: **for**  $c \leftarrow cost(\epsilon) + 1$  **to**  $maxCost$  **do**  
8:      $questions =$  buildQuestionMark( $c - cost(?)$ )  
9:      $stars =$  buildStar( $c - cost(*)$ )  
10:      $concats =$  buildConcat( $c - cost(\cdot)$ )  
11:      $unions =$  buildUnion( $c - cost(+)$ )  
12:      $languageCache[c] =$  questions ++ stars ++ concats ++ unions ▷ ++ is concatenation  
13: **return** "not\_found" ▷ Procedures in loop will return solution directly to caller of main, if found

---

# Language Cache

We construct all needed languages bottom-up, from lower to higher cost, and keep the constructed languages in memory for later re-use:



---

**Algorithm 2** Pseudocode for concatenation (*buildConcat* procedure in Alg. 1)

---

```
1: Input cost  $c$ , globals used: languageCache,  $P$ ,  $N$ 
2: Output A list of new CSs generated by concatenation
3: 

---


4: outList  $\leftarrow []$ 
5: for all  $L, R$  such that  $L + R = c$  do
6:   for all  $lCS \in \text{languageCache}(L)$  do
7:     for all  $rCS \in \text{languageCache}(R)$  do
8:        $i \leftarrow 1$ 
9:        $newCS \leftarrow 0$ 
10:      for  $w \leftarrow 0$  to  $\#ic(P \cup N) - 1$  do
11:        for all pair  $(l, r) \in \text{gt}[w]$  do
12:          if  $(lCS \& l) \neq 0$  and  $(rCS \& r) \neq 0$  then
13:             $newCS \leftarrow newCS | i$ 
14:             $i \leftarrow i \ll 1$ 
15:             $isUnique \leftarrow \text{hashSet.insert}(newCS)$ 
16:            if  $isUnique$  then
17:              if  $newCS \models (P, N)$  then
18:                print  $newCS$  and terminate program
19:              outList.insert(newCS)
20: return outList
```

---

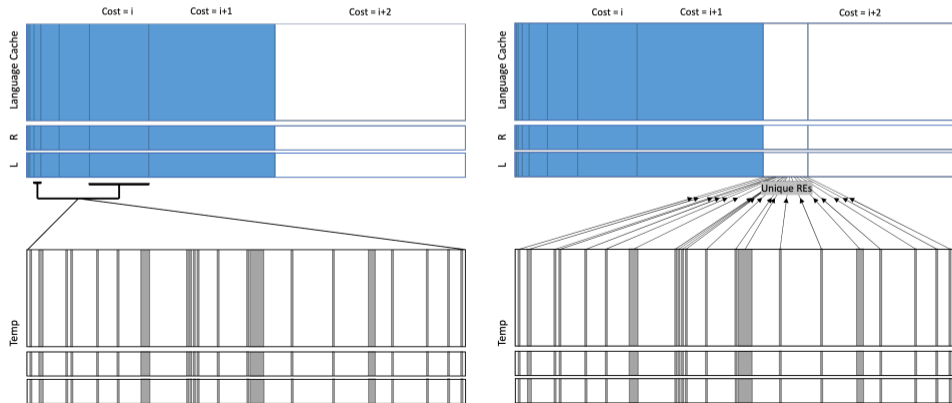
All local, except global uniqueness check



In order to maximise parallelism on GPUs, we allocate each new characteristic sequence (which may be redundant) into a temporary array

Each new characteristic sequence is composed from older entries in the language cache by simple, fast and **local** bitvector operations

For **each** new characteristic sequence we run a **global** uniqueness check



# Measurements

AlphaRegex benchmarks (previously state-of-the-art) are too small.

Our new, large benchmarks suite

- Alphabet  $\Sigma$ ,
- $le$  is the maximal length of example strings,
- $p$  and  $n$ , the numbers of positive and negative examples, respectively.

With those parameters, we define two complementary benchmark generation schemes. Both create instances  $(P, N)$  by sampling uniformly from two different spaces of random strings.

- TYPE 1:  $\{(P, N) \in \Sigma^{\leq le} \times \Sigma^{\leq le} \mid \forall w \in P \cup N. \#P = p, \#N = n, P \cap N = \emptyset\}$
- TYPE 2:  $\{((P_0, \dots, P_{le}), (N_0, \dots, N_{le})) \in Y \times Y \mid \Sigma_i \#P_i = p, \Sigma_i \#N_i = n, \forall i. P_i \cap N_i = \emptyset\}$

Input					CPU	GPU			
Type	No	# P	# N	Cost Function	Sec	Sec	Speed-up	# REs	
1	50	10	12	(1, 1, 1, 1, 1)	5080.7850	4.9512	<b>1026x</b>	26,774,099,142	
1	51	12	9	(10, 1, 1, 1, 1)	4699.8137	4.4966	<b>1045x</b>	23,824,118,297	
1	73	10	11	(1, 10, 1, 1, 1)	5805.2168	3.7144	<b>1562x</b>	22,703,639,676	
1	20	9	9	(1, 1, 10, 1, 1)	2893.4835	2.8935	<b>1000x</b>	13,567,472,188	
1	73	10	11	(1, 1, 1, 10, 1)	2901.9297	2.9504	<b>983x</b>	11,706,686,339	
1	31	8	9	(1, 1, 1, 1, 10)	5856.6925	3.9973	<b>1465x</b>	14,210,157,835	
1	57	12	10	(10, 10, 10, 10, 1)	2804.6793	3.4322	<b>817x</b>	14,163,906,090	
1	50	10	12	(10, 10, 10, 1, 10)	4519.9456	4.9096	<b>920x</b>	23,349,552,935	
1	57	12	10	(10, 10, 1, 10, 10)	4301.8548	4.5243	<b>950x</b>	20,257,045,497	
1	97	12	12	(10, 1, 10, 10, 10)	5608.7286	4.7782	<b>1173x</b>	19,680,542,658	
1	61	12	10	(1, 10, 10, 10, 10)	2915.0938	3.0532	<b>954x</b>	14,322,039,866	
1	88	12	9	(20, 20, 20, 5, 30)	6899.0045	4.6904	<b>1470x</b>	25,193,577,825	
2	88	14	8	(1, 1, 1, 1, 1)	3783.9772	4.2462	<b>891x</b>	23,697,549,545	
2	150	14	12	(10, 1, 1, 1, 1)	4228.2773	4.4120	<b>958x</b>	23,125,803,623	
2	158	12	14	(1, 10, 1, 1, 1)	2975.9956	2.4887	<b>1195x</b>	11,432,891,412	
2	136	11	14	(1, 1, 10, 1, 1)	3374.8873	3.6080	<b>935x</b>	18,241,755,827	
2	107	12	12	(1, 1, 1, 10, 1)	2432.4320	4.4120	<b>551x</b>	24,954,272,802	
2	32	10	7	(1, 1, 1, 1, 10)	7400.8135	4.6482	<b>1592x</b>	16,729,795,052	
2	136	11	14	(10, 10, 10, 10, 1)	2907.9182	3.9689	<b>732x</b>	17,476,988,322	
2	200	13	8	(10, 10, 10, 1, 10)	9687.7952	4.5366	<b>2135x</b>	6,037,014,423	
2	107	12	12	(10, 10, 1, 10, 10)	3383.1937	4.5071	<b>750x</b>	20,697,274,025	
2	81	8	14	(10, 1, 10, 10, 10)	3497.9013	4.6699	<b>749x</b>	21,869,903,022	
2	88	14	8	(1, 10, 10, 10, 10)	3405.5536	4.1602	<b>818x</b>	21,889,508,744	
2	158	12	14	(20, 20, 20, 5, 30)	5804.8112	4.9228	<b>1179x</b>	23,163,079,580	
<b>Average</b>					4465.4493	4.1238	<b>1077x</b>	19,127,861,447	

5160 benchmarks (200 Type 1, 230 Type 2, 12 cost functions). Examples hardest within timeout < 5 sec on GPU. GPU: Nvidia A100-SXM4, CPU: Intel Xeon, 2.20 GHz, 25 GB

## Future work

Algorithmic engineering, is uniqueness check for every new characteristic sequence optimal?

Context-free inference? (implemented on CPU, in progress: GPU)

Scaling: approximate REI (in progress)

Mathematics of synthesis, e.g.

$$\frac{\text{Regular grammar}}{\text{Infix-closure}} = \frac{\text{General grammar}}{\text{???-closure}}$$

Comparison with LLMs (paper under anonymous review)

# Thanks

Code: <https://github.com/MojtabaValizadeh/paresy>

Paper: <https://arxiv.org/abs/2305.18575>

Code:



Paper:

