



# Scalable Synthesis of Regular Expressions from Only Positive Examples

`a{5}[a-z]{2,}`

`(+)?\d+`

## Only Positive Examples

Mark Barbone\* and Elizaveta Pertseva\*  
Advisors: Nadia Polikarpova and Taylor Berg-Kirkpatrick

`\d{2}.\d{2}.\d{4}`

`[a-z]+`

### Defining the Problem

#### Motivating Example



- Alex is searching for Brazilian CNPJs in a large file
- Quickly finds a couple examples
- What now? Regex+

`02.916.265/0001-60`

`60.701.190/0001-04`

Regex+

`\d{2}\.\d{3}\.\d{3}/0001-\d{2}`

#### Challenges

1. Underspecification

“Correct” but too **specific**:  
`(02\.916\.265/0001-60)` OR  
`(60\.701\.190/0001-04)`

“Correct” but too **simple**: `.*`

2. Search

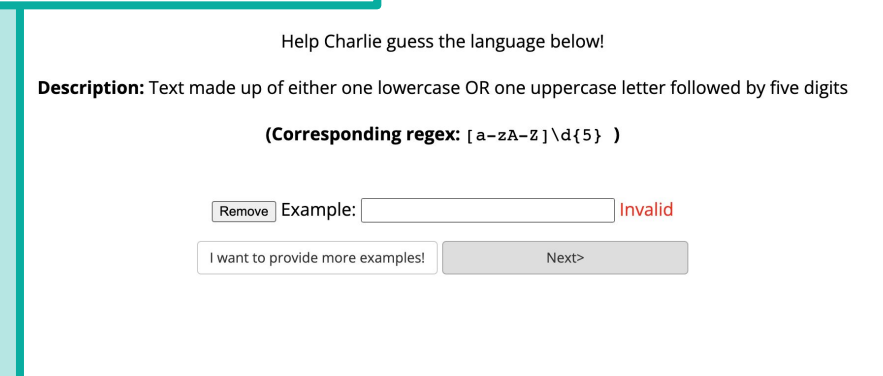
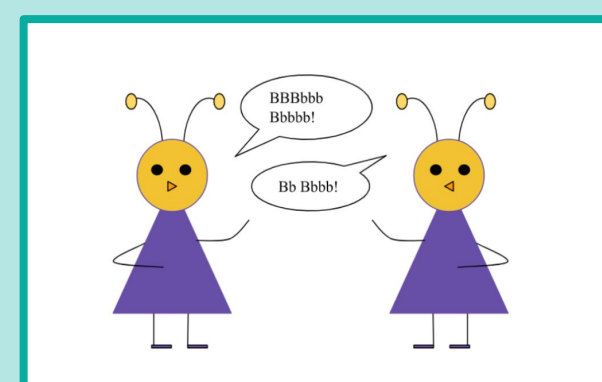
### Contribution 1: Pragmatic Ranking Function

#### Model

$$P(\text{regex}|\text{input}) \sim P(\text{input}|\text{regex}) \cdot P(\text{regex})$$

specificity

simplicity

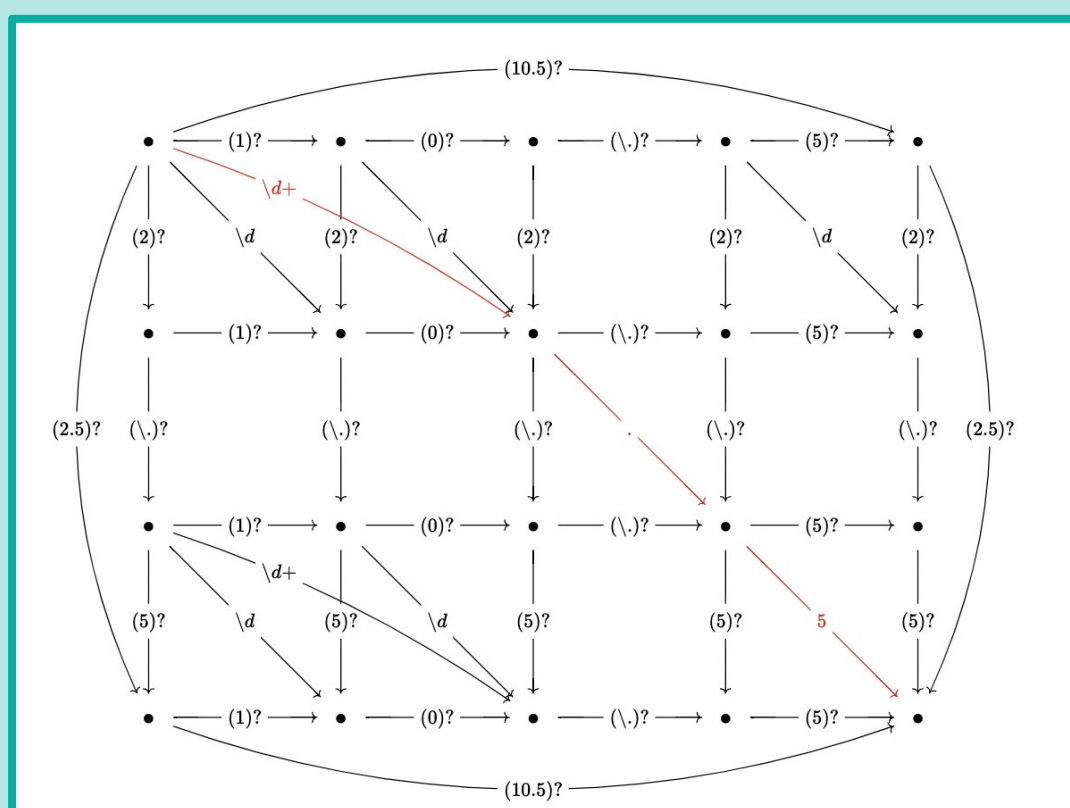


#### Studying Reality

- Human study to evaluate hypothesis posed as a game
- 412** new data points

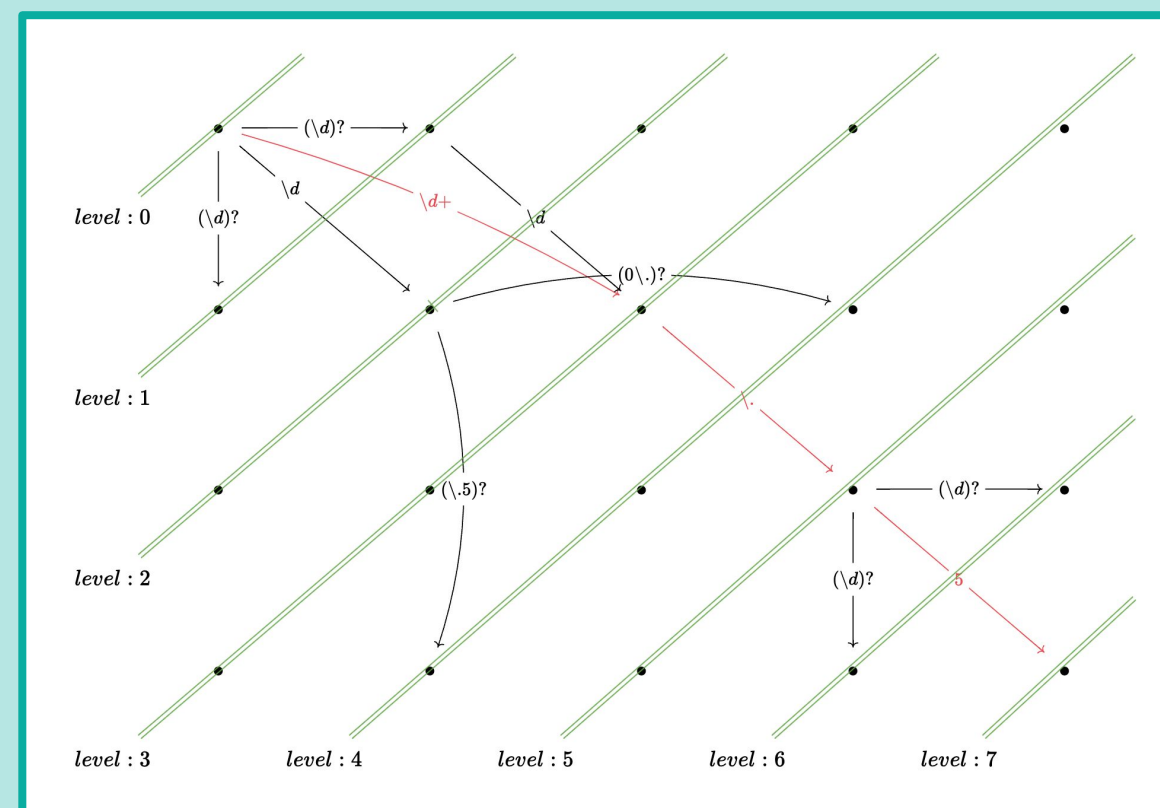
### Contribution 2: Search Algorithms

#### VSA



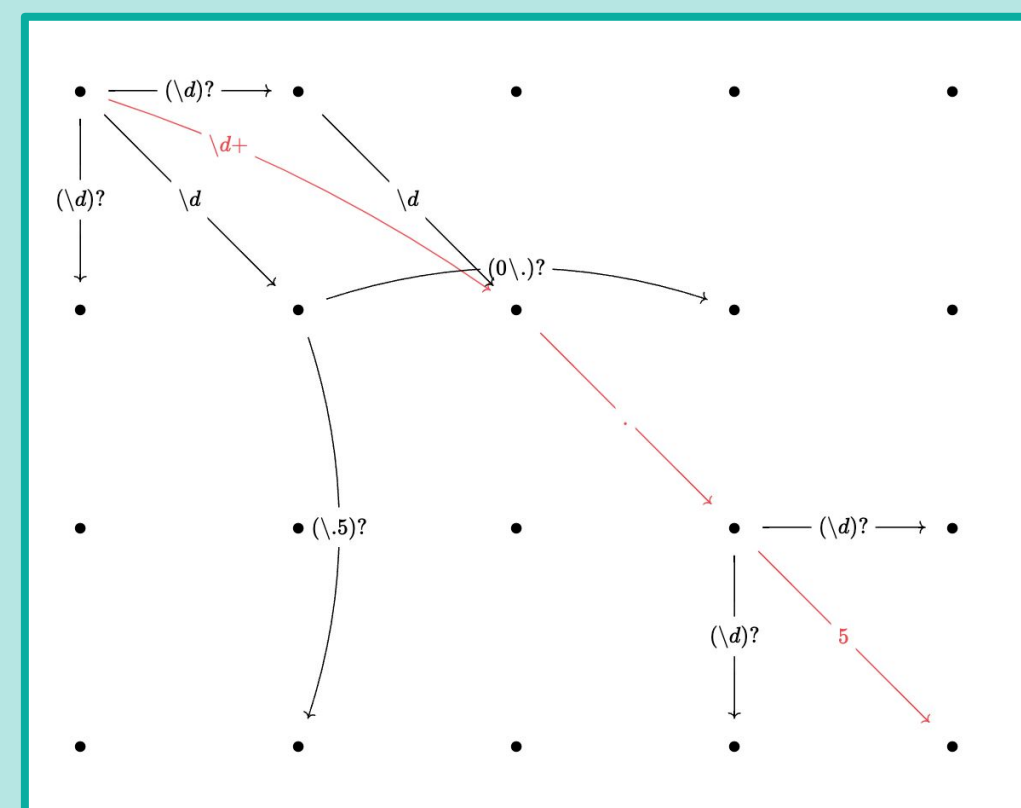
Slow, memory intensive **✗**  
Completeness Guarantees **✓**

#### Beam



Fast **✓**  
Completeness Guarantees **✗**

#### A\*



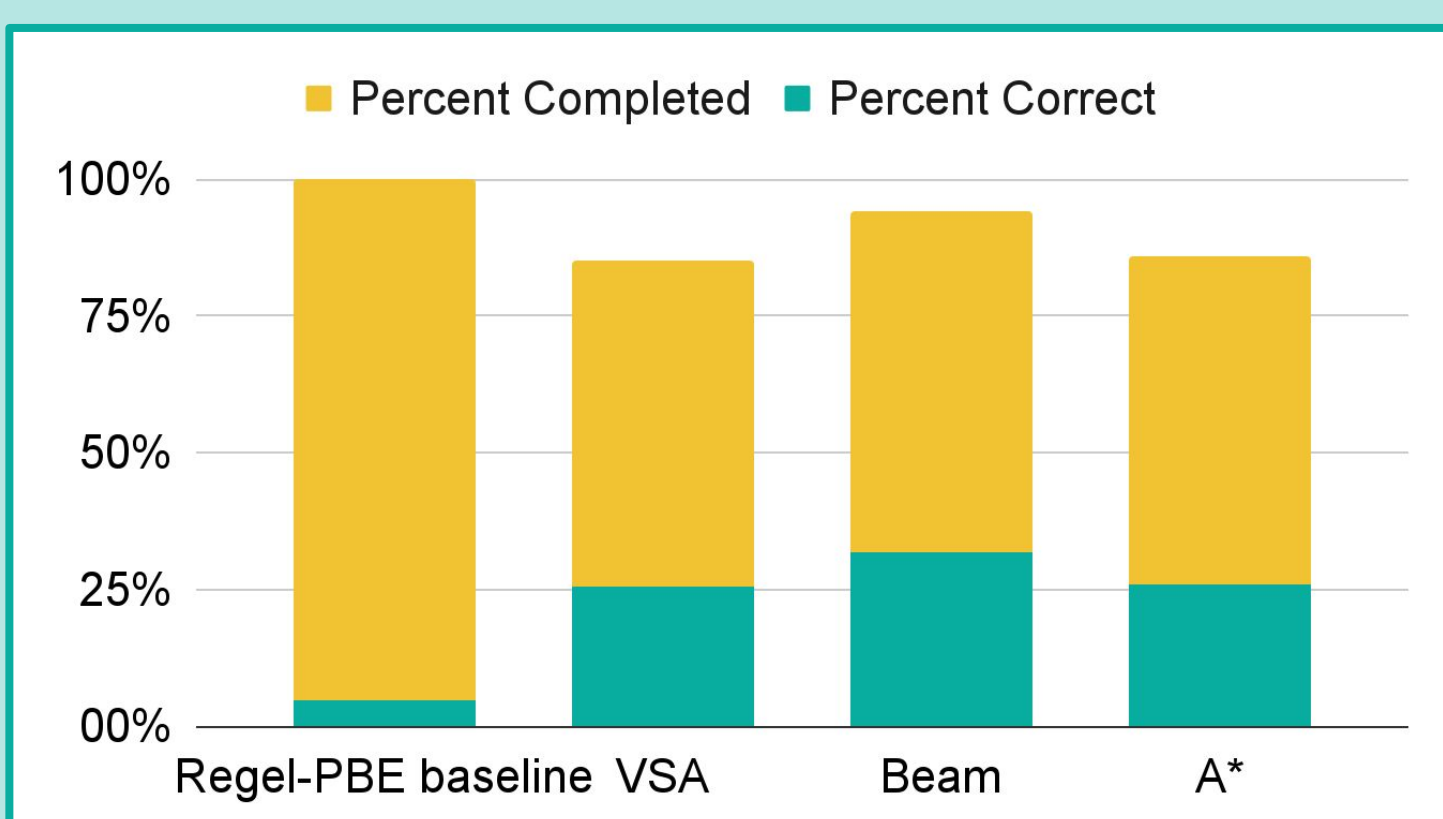
Moderately Fast **✓**  
Completeness Guarantees **✓**

admissible heuristic:

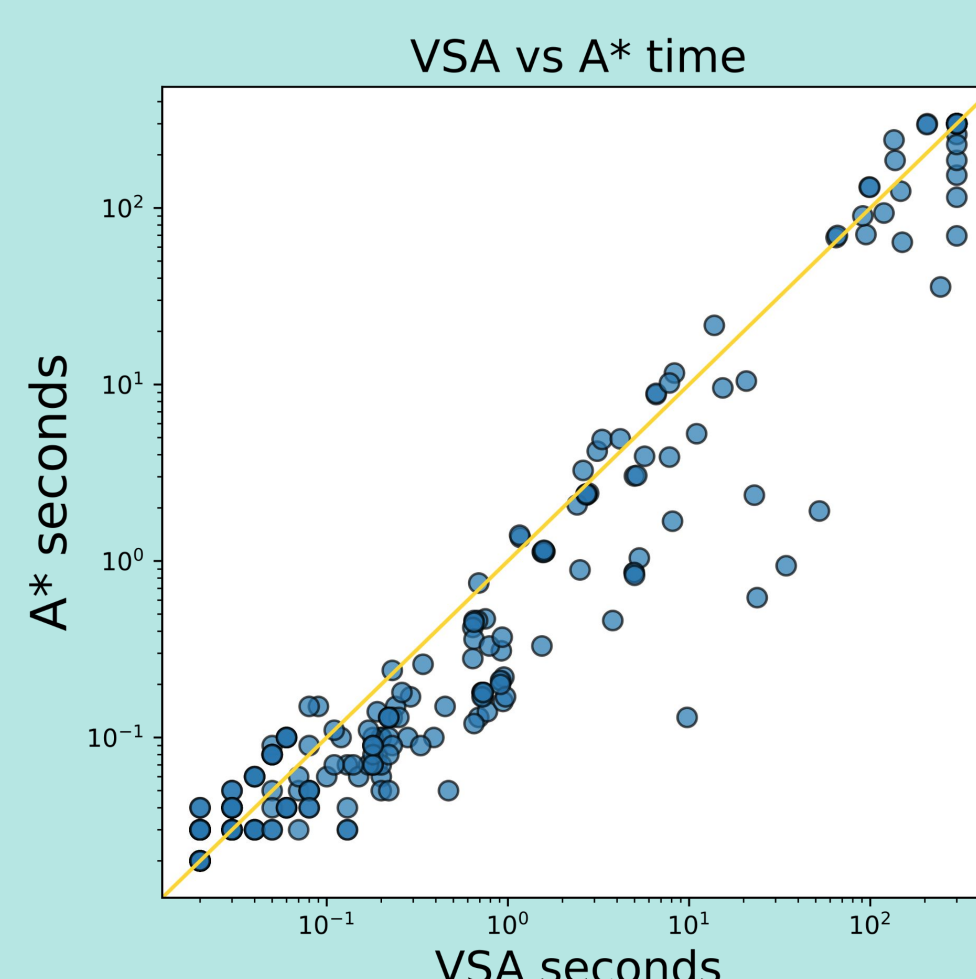
$$\min_{\text{regexes } R} \left[ \text{SIMPL}(R) + \sum_{i=1}^N \text{SPEC}(R, e_i) \right]$$

$$\geq \max_{\text{examples } e_i} \left[ \min_{\text{regexes } R} \left[ \text{SIMPL}(R) + \text{SPEC}(R, e_i) \right] \right]$$

### Comparison



### Evaluation



### Specific Example

111-324-2344  
276-943-3044  
584-143-2455  
000-000-0000  
249-394-1232

Correct	<code>\d{3}-\d{3}-\d{4}</code>
Regal-PBE	<code>([0-9]).*</code>
VSA	MEMORY LIMIT
A*	<code>\d{3}-\d{3}-\d{4}</code>
Beam	<code>\d{3}-\d{3}-\d{4}</code>