

# Synthesizing Extraction Rules from User Examples with SEER



Maeda F. Hanafi (maeda.hanafi@nyu.edu), Azza Abouzied (azza@nyu.edu),  
Laura Chiticariu (chiti@us.ibm.com), Yunyaoli Li (yunyaoli@us.ibm.com)

## Challenges in Information Extraction (IE)

- **IE is a key technology for understanding text**
  - Companies extract data to build knowledge-bases.
  - Journalists extract crime rates from FBI datasets to analyze changes.
  - Business analysts extract people and corporations from SEC filings to analyze relationships.

- **Two approaches to IE:**

1. Machine learning:
  - Cons: Requires large training datasets.
  - Cons: Produces hard-to-interpret statistical models.
  - Pros: Once training data is given, quick to generate results.
2. Programming:
  - Cons: Requires learning the language.
  - Cons: Labor-intensive and time-consuming.
  - Pros: Transparency

## Solution: SEER

- **SEER: Combines the best of both approaches.**
  - Learning high-quality extraction rules from handful of examples.
- **SEER synthesizes and suggests visual rules to the user.**
  - Quick to get high-quality extraction rules and results.
  - Reduced manual development effort.
- Extraction Rules = sequences of *primitives*  
Primitives capture certain tokens.

**Pre-builts** extract entities: P: City

**Literals** extract exact strings: L: 'Dubai'

**Dictionaries** extract sets of exact strings: D: {Dubai, London}

**Regexes** extract from a library of regexes: R: [A-Za-z]+

**Token gaps** extract any token: T: 0-1

## 1. Highlight and Process Examples

- **The user highlights examples of text to extract.**

Law enforcement agencies submitted incident reports involving 5,479 criminal incidents. In cities with populations from 50,000 to 99,000 inhabitants, violent crimes **rose 11 percent**. In cities of populations above 100 thousand, hate crime offenses **decreased 5 percent in 2010 alone**.

Positive examples:    
Negative examples:  

- **SEER enumerates primitives for each positive token.**

Tokens:	Primitives:	Tokens:	Primitives:
rose	<span>L: 'rose'</span> <span>R: [A-Za-z]+</span>	decreased	<span>L: 'decreased'</span> <span>R: [A-Za-z]+</span>
11	<span>P: Number</span> <span>P: Integer</span> <span>L: '11'</span> <span>R: [0-9]+</span> <span>T: 0-1</span>	5	<span>P: Number</span> <span>P: Integer</span> <span>L: '5'</span> <span>R: [0-9]+</span> <span>T: 0-1</span>
percent	<span>L: 'percent'</span> <span>R: [A-Za-z]+</span>	percent	<span>L: 'percent'</span> <span>R: [A-Za-z]+</span>
11 percent	<span>P: Percentage</span>	5 percent	<span>P: Percentage</span>

- **Assign scores to the primitives to guide the search.**
  - Intuition: Rule developers prefer certain primitives for certain tokens.
  - If the primitive's tokens have inherent semantics:

Semantic Token: Dubai

Primitive types:	Score:	Primitives:
Prebuilt	1.0 to 0.6	<span>P: City = 1.0</span>
Dictionary, Literal	0.4	<span>L: 'Dubai' = 0.4</span>
Regex, Token Gap	0.2	<span>R: [A-Za-z]+ = 0.2</span>

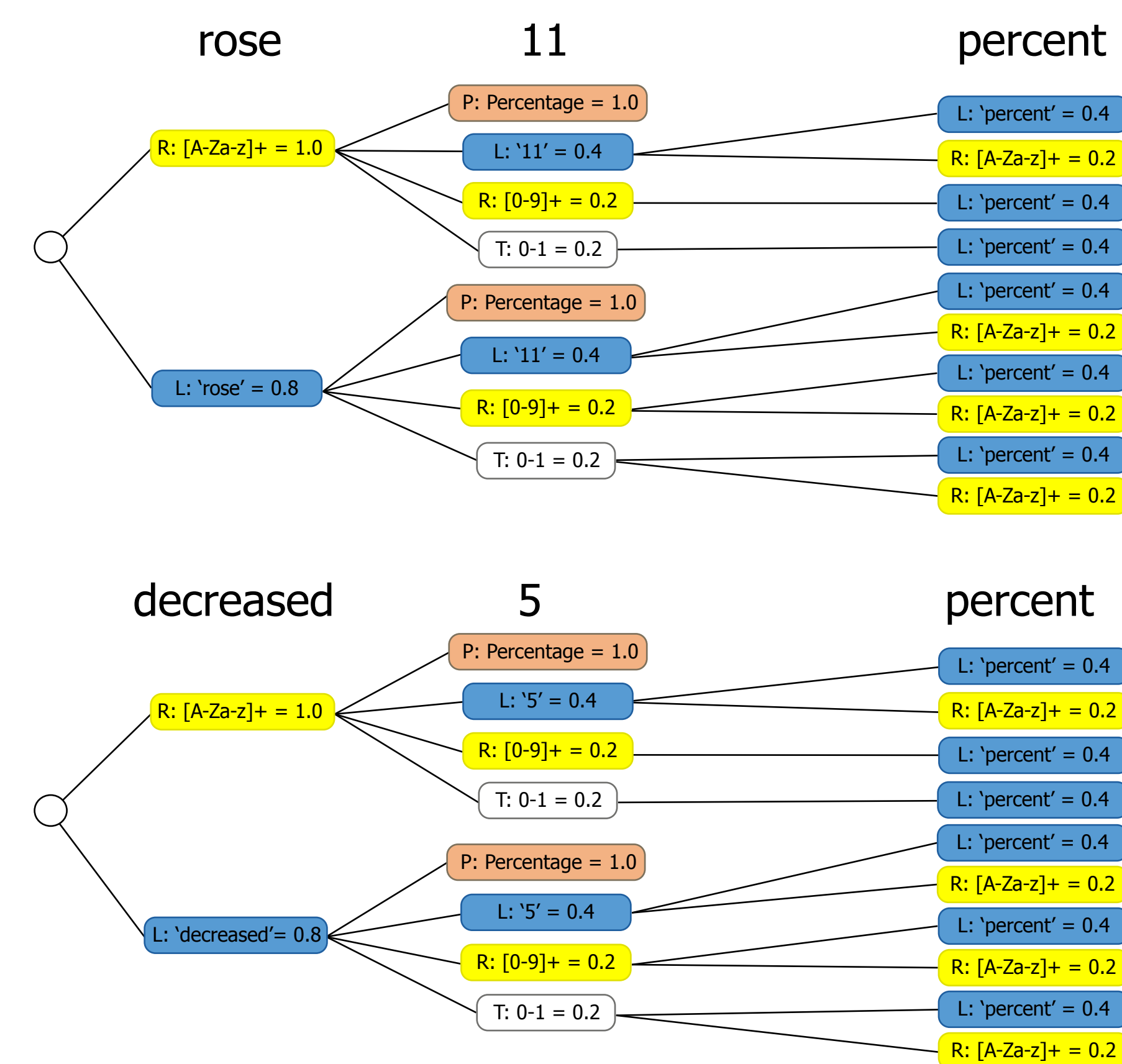
- If the primitive's tokens are syntactic:

Syntactic Token: of

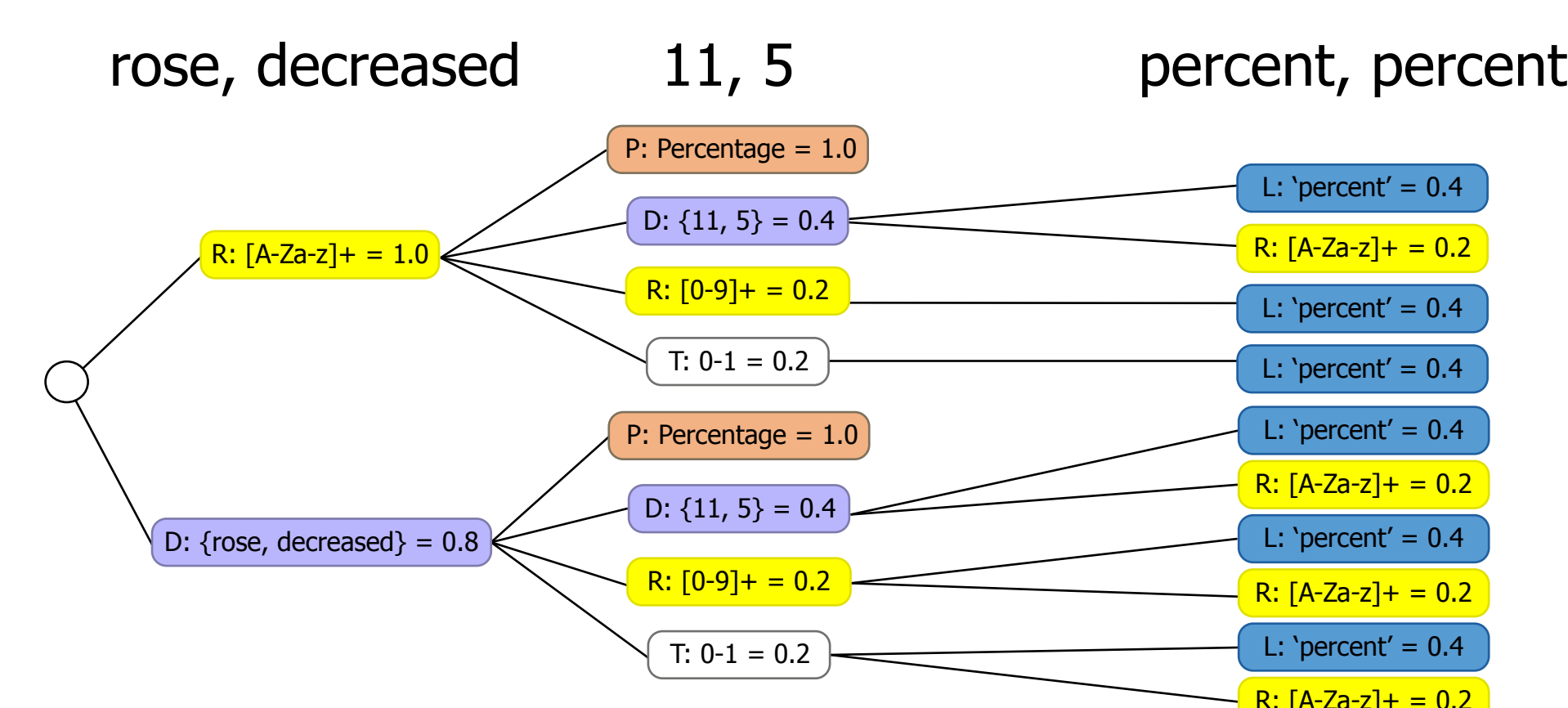
Primitive types:	Score:	Primitives:
Regex, Token Gap	1.0	<span>R: [A-Za-z]+ = 1.0</span>
Dictionary, Literal	0.8	<span>L: 'of' = 0.8</span>

## 2. Synthesize Rules

- **For each positive example, generate a tree of rules.**
  - Extraction Rule = path of primitives from root to leaf
- Rules capturing negative examples are removed.

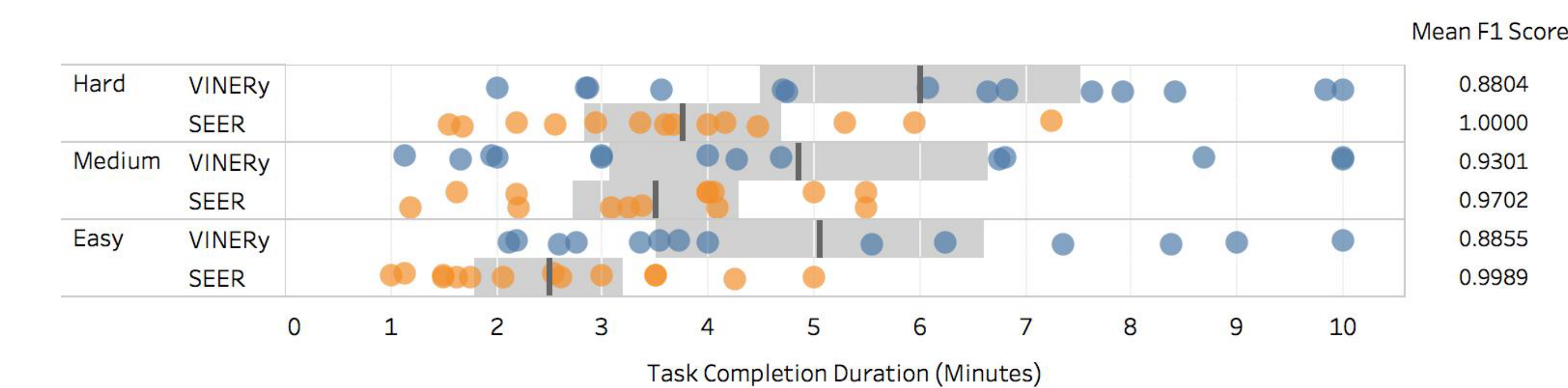


- **Capture all user examples by intersecting trees.**
  - Generalize the rules to capture all user examples.



## 4. Evaluations

- **Baseline: VINERy Visual Programming IDE for IE** (VLDB '15)
- User Study: Each user used both tools on two datasets in random order to minimize learning effects.
- Black bars = mean duration.  
Gray bars = 95% mean confidence interval.



## 3. Prune and Suggest Rules

- **Group similar rules to ensure rule diversity.**
  - Similar rules have the same primitive types.
  - Pick the top rule of each group.
- Rule score = AVG(Primitive Scores)

Groups:

Regex, Pre-built	<span>R: [A-Za-z]+ = 1.0</span> <span>P: Percentage = 1.0</span> = 1.00
Regex, Dictionary, Literal	<span>R: [A-Za-z]+ = 1.0</span> <span>D: {11, 5} = 0.4</span> <span>L: 'percent' = 0.4</span> = 0.60 <span>D: {rose, decreased} = 0.8</span> <span>R: [0-9]+ = 0.2</span> <span>L: 'percent' = 0.4</span> = 0.47
Regex, Dictionary	<span>R: [A-Za-z]+ = 1.0</span> <span>D: {11, 5} = 0.4</span> <span>R: [A-Za-z]+ = 0.2</span> = 0.53 <span>D: {rose, decreased} = 0.8</span> <span>D: {11, 5} = 0.4</span> <span>R: [A-Za-z]+ = 0.2</span> = 0.47 <span>D: {rose, decreased} = 0.8</span> <span>R: [0-9]+ = 0.2</span> <span>R: [A-Za-z]+ = 0.2</span> = 0.40
Regex, Literal	<span>R: [A-Za-z]+ = 1.0</span> <span>R: [0-9]+ = 0.2</span> <span>L: 'percent' = 0.4</span> = 0.53
Regex, Token Gap, Literal	<span>R: [A-Za-z]+ = 1.0</span> <span>T: 0-1 = 0.2</span> <span>L: 'percent' = 0.4</span> = 0.53
Dictionary, Pre-built	<span>D: {rose, decreased} = 0.8</span> <span>P: Percentage = 1.0</span> = 0.90
Dictionary, Literal	<span>D: {rose, decreased} = 0.8</span> <span>D: {11, 5} = 0.4</span> <span>L: 'percent' = 0.4</span> = 0.53
Dictionary, Token Gap, Literal	<span>D: {rose, decreased} = 0.8</span> <span>T: 0-1 = 0.2</span> <span>L: 'percent' = 0.4</span> = 0.47
Dictionary, Token Gap, Regex	<span>D: {rose, decreased} = 0.8</span> <span>T: 0-1 = 0.2</span> <span>R: [A-Za-z]+ = 0.2</span> = 0.40

- **Show the top ten rules to the user.**

1. R: [A-Za-z]+ = 1.0 P: Percentage = 1.0 = 1.00
2. D: {rose, decreased} = 0.8 P: Percentage = 1.0 = 0.90
3. R: [A-Za-z]+ = 1.0 D: {11, 5} = 0.4 L: 'percent' = 0.4 = 0.60
4. R: [A-Za-z]+ = 1.0 D: {11, 5} = 0.4 R: [A-Za-z]+ = 0.2 = 0.53
5. R: [A-Za-z]+ = 1.0 R: [0-9]+ = 0.2 L: 'percent' = 0.4 = 0.53
6. R: [A-Za-z]+ = 1.0 T: 0-1 = 0.2 L: 'percent' = 0.4 = 0.53
7. D: {rose, decreased} = 0.8 D: {11, 5} = 0.4 L: 'percent' = 0.4 = 0.53
8. D: {rose, decreased} = 0.8 T: 0-1 = 0.2 L: 'percent' = 0.4 = 0.47
9. D: {rose, decreased} = 0.8 T: 0-1 = 0.2 R: [A-Za-z]+ = 0.2 = 0.40