



RegexScalpel: Regular Expression Denial of Service (ReDoS) Defense by Localize-and-Fix

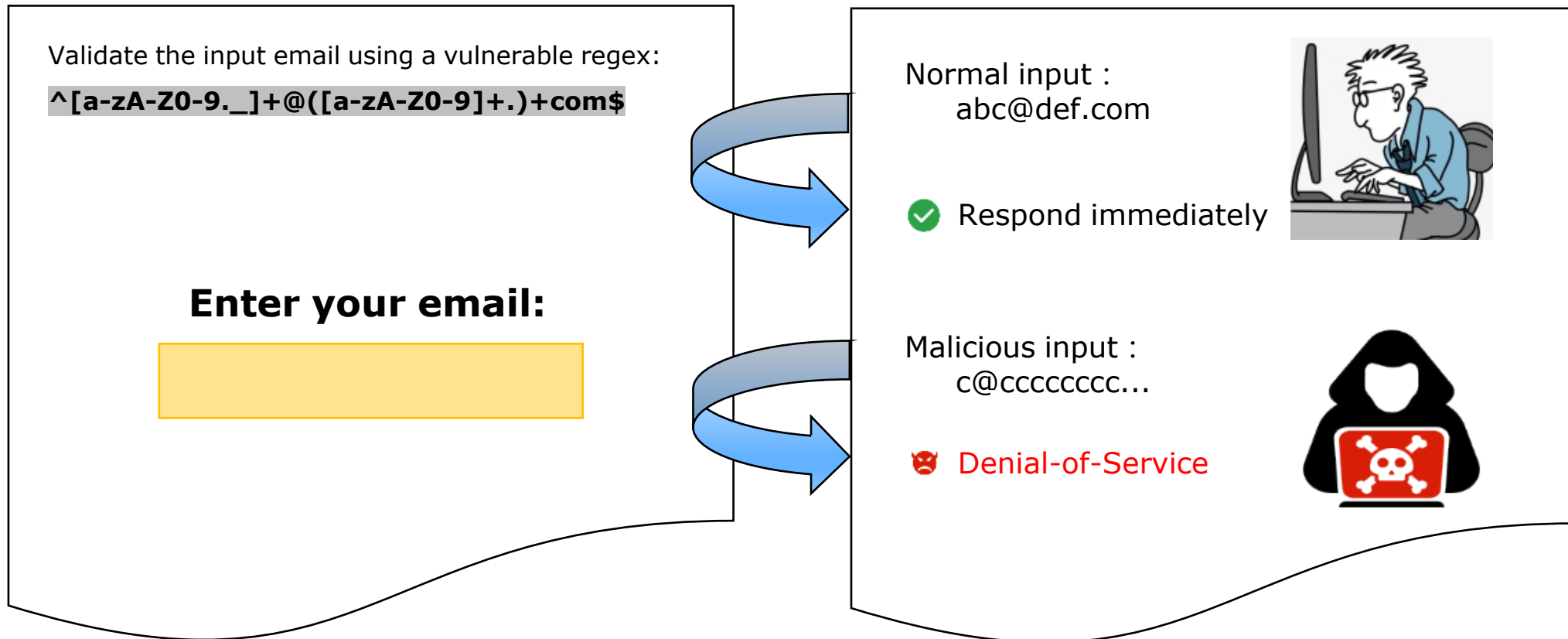
Yeting Li, **Yecheng Sun**, Zhiwu Xu, Jialun Cao, Yuekang Li, Rongchen Li,
Haiming Chen, Shing-Chi Cheung, Yang Liu, Yang Xiao

Institute of Information Engineering, Chinese Academy of Sciences
Institute of Software, Chinese Academy of Sciences
University of Chinese Academy of Sciences
Shenzhen University
The Hong Kong University of Science and Technology
Nanyang Technological University

Regular Expression Denial-of-Service

The Server Side

The Client Side



Regular Expression Denial-of-Service Defense

- Existing Solutions

- **Regex engine substitution**

Omitting extended features, consuming space, bringing semantic differences or incompatibilities.

- **Input length restriction**

Facing a dilemma known as "Goldilocks problem".

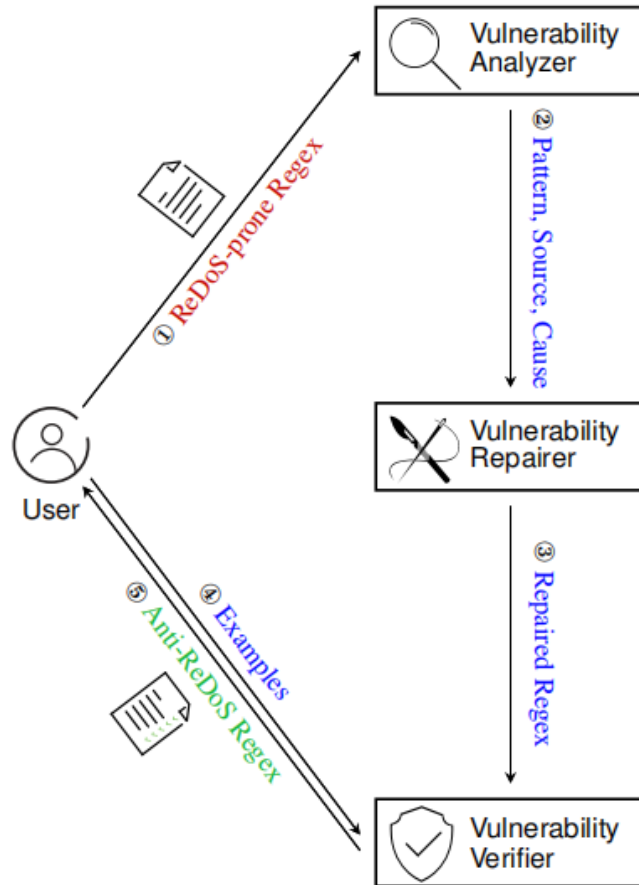
- **Regex repair**

- Repairing vulnerable regexes can greatly mitigate their vulnerabilities.
 - But it's challenging.

Table 13: The proportion of Maintainers' Defense Actions.

Defense Strategy	SOLA-DA	CVE	Total
Regex Engine Substitution	0 (0%)	1 (0.24%)	1 (0.22%)
Input Length Limit	1 (2.94%)	2 (0.48%)	3 (0.67%)
Code Logic Modification	6 (17.65%)	17 (4.11%)	23 (5.13%)
Regex Repair	21 (61.76%)	392 (94.69%)	413 (92.19%)
No Fix	6 (17.65%)	2 (0.48%)	8 (1.79%)
#Regex	34	414	448

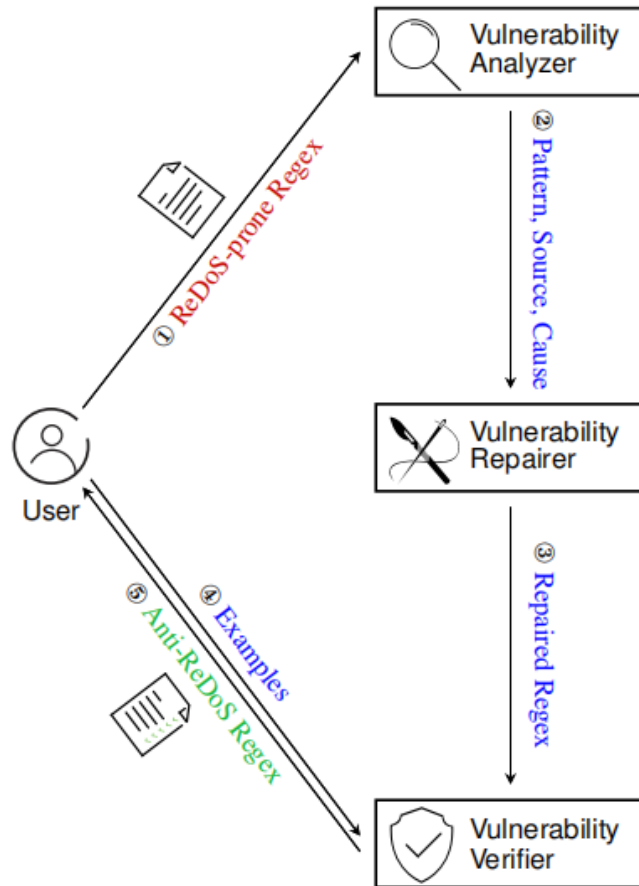
Regular Expression Denial-of-Service Defense



- Our Solution

- We proposed RegexScalpel, a regex ReDoS vulnerability analysis and repair framework based on localize-and-fix.
- RegexScalpel can preserve the semantics of the original regex, and the iterative repair method also keeps out vulnerabilities of the repaired regexes.

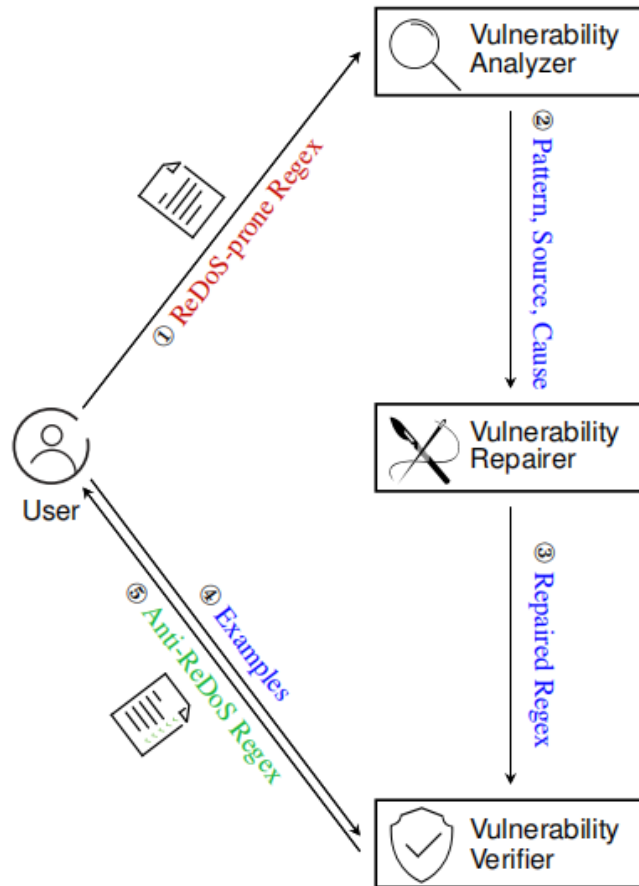
Regular Expression Denial-of-Service Defense



- Our Solution

- RegexScalpel first leverages the **fine-grained vulnerability patterns** to localize the vulnerabilities and obtain the information for the repair.
- The information includes their **vulnerable patterns**, the **source** (i.e., the pathological sub-regexes), and the **root causes** (e.g., the overlapping sub-regexes).

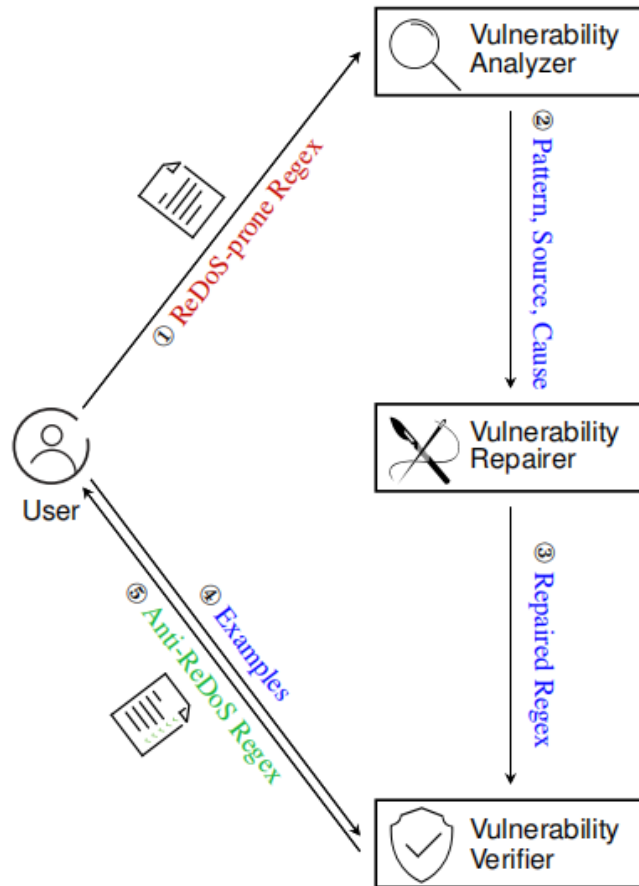
Regular Expression Denial-of-Service Defense



- Our Solution

- RegexScalpel then fixes the pathological sub-regexes according to the **repair patterns** and the information.
- The repair patterns use **micro-manipulations** (e.g., adding a lookahead, deleting a quantifier or sub-regex, modifying a quantifier or sub-regex) to remove the overlapping paths or reducing the maximum times of backtracking.

Regular Expression Denial-of-Service Defense



- Our Solution

- RegexScalpel next determines whether the repaired regexes are ReDoS-invulnerable and whether it can pass all the given test cases.
- If so, the repaired regex is called a successful one. If not, RegexScalpel continues the vulnerability analysis and repairs it.
- RegexScalpel finally returns a repaired regex randomly chosen from the successful ones.

Vulnerable Patterns

Table 1: The Sub-pattern, Vulnerability Description, Example Regex, and Results from AnaNQ of the Pattern $\mathcal{N}Q$.

No.	Sub-pattern	Vulnerability Description	Example Regex	Returned Triple
#1	$\mathcal{N}Q_1$	$r = r_1\{m,n\}$, where $r_1 = r_p\{m_p, n_p\}$, $n_p > 1$, and $n > 1$	$\delta_1 = \backslash[(\backslash d+;)?(\backslash d+)*m$ (CVE-2015-9239)	$(\mathcal{N}Q_1, (\backslash d+)^*, [+,*])$
#2	$\mathcal{N}Q_2$	$r = (r_0r_1r_2)\{m,n\}$, where r_0 and r_2 are nullable, $r_1 = r_p\{m_p, n_p\}$, $n_p > 1$, and $n > 1$	$\delta_2 = ^(\text{https? ftp}):\backslash\backslash(-\backslash.)?([\backslash s\backslash/?\backslash.\#]+ \backslash.?) + (\backslash/[\backslash s]*)?[\backslash s\backslash.,,]\$$ (CVE-2021-26272)	$(\mathcal{N}Q_2, ([\backslash s\backslash/?\backslash.\#]+\backslash.?)^+, [+,*])$
#3	$\mathcal{N}Q_3$	$r = (r_0r_1r_2)\{m,n\}$, where r_0 and r_2 are nullable, $r_1 = (\dots r_p\{m_p, n_p\} \dots)$, $n_p > 1$, and $n > 1$	$\delta_3 = D[oD]?(\backslash[[\backslash \backslash \backslash]]*\backslash \backslash s+)+MMMM?$ (moment)	$(\mathcal{N}Q_3, (\backslash[[\backslash \backslash \backslash]]*\backslash \backslash s+)^+, [+,*])$

- Nested Quantifiers (NQ)

- The NQ pattern is a regex with nested quantifiers.
- In order to facilitate fixing the pathological regex, we subdivided NQ pattern into three sub-patterns (i.e., NQ1, NQ2 and NQ3).

Vulnerable Patterns

Table 2: The Sub-pattern, Vulnerability Description, Example Regex, and Results from AnaQOD of the Pattern QOD .

No.	Sub-pattern	Vulnerability Description	Example Regex	Returned Triple
#1	QOD_1	$r = \alpha\{m, n\}$, $\alpha = (r_1 \dots r_k)$, and $\exists \mathcal{L}(\alpha_1) \cap \mathcal{L}(\alpha_2) \neq \emptyset$, where $\alpha_1 = r_{p_1} r_{p_2} \dots r_{p_t}$, $\alpha_2 = r_{q_1} r_{q_2} \dots r_{q_s}$, $1 \leq t, s \leq k$, $1 \leq i \leq t$, $1 \leq j \leq s$, $1 \leq p_i, q_j \leq k$, and $p_i \neq q_j$	$\delta_4 = ^{(\backslash d \backslash . \backslash . \backslash d \backslash d)} + \$$ ([14])	$(QOD_1, (\backslash d \backslash . \backslash . \backslash d \backslash d) +, [\backslash d \backslash . \backslash d, \backslash d \backslash . \backslash d])$
#2	QOD_2	$r = \alpha\{m, n\}$, $\alpha = (r_1 \dots r_k)$, and $\exists \mathcal{L}(r_p) \cap \mathcal{L}(\alpha_1) \neq \emptyset$, $\mathcal{L}(\eta) \cap \mathcal{L}(r_p) = \emptyset$, where $\alpha_1 = r_{q_1} r_{q_2} \dots r_{q_s}$, $1 \leq i \leq k$, $r_i = r_l r_p$, r_l is not nullable, $1 \leq s \leq k$, $1 \leq j \leq s$, $1 \leq q_j \leq k$	$\delta_5 = (?: [\backslash w -] \backslash \$ [- \backslash w] + \# \backslash \{ \backslash \$ [- \backslash w] + \backslash \}) + (? = \backslash s * :)$ (CVE-2021-23341)	$(QOD_2, (?: [\backslash w -] \backslash \$ [- \backslash w] + \# \backslash \{ \backslash \$ [- \backslash w] + \backslash \}) +, [[\backslash w -], \backslash \$ [- \backslash w] +])$

- Quantified Overlapping Disjunction (QOD)
 - The QOD pattern is a quantified disjunction whose multiple inner sub-regexes overlap.
 - In order to facilitate fixing the pathological regex, we subdivided QOD pattern into two sub-patterns (i.e., QOD1 and QOD2).

Vulnerable Patterns

Table 3: The Sub-pattern, Vulnerability Description, Example Regex, and Results from AnaQOA of the Pattern QOA .

No.	Sub-pattern	Vulnerability Description	Example Regex	Returned Triple
#1	QOA_1	$r = (...r_1r_2...)\{m, n\}$, where $r_1 = r_p\{m_p, n_p\}$, $r_2 = r_q\{m_q, n_q\}$, $\mathcal{L}(r_1) \cap \mathcal{L}(r_2) \neq \emptyset$	$\delta_6 = \text{^_}[\text{^\wedge}\backslash\text{w_}] + \backslash\text{w} + \text{_}\$$ (pylint)	$(QOA_1, [\text{^\wedge}\backslash\text{w_}] + \backslash\text{w} +, [[\text{^\wedge}\backslash\text{w_}] +, \backslash\text{w} +])$
#2	QOA_2	$r = (...r_1r_2r_3...)\{m, n\}$, where $r_2 = r_t\{m_t, n_t\}$ is nullable, $r_1 = r_p\{m_p, n_p\}$, $r_3 = r_q\{m_q, n_q\}$, $\mathcal{L}(r_1) \cap \mathcal{L}(r_3) \neq \emptyset$	$\delta_7 = \text{^}(>=? <=?)\backslash\text{s}^*(\backslash\text{d}^*\backslash\text{.}? \backslash\text{d} +) \text{\%}\$$ (CVE-2021-23364)	$(QOA_2, \backslash\text{d}^*\backslash\text{.}? \backslash\text{d} +, [\backslash\text{d}^*, \backslash\text{d} +])$
#3	QOA_3	$r = (r_1...r_2)\{m, n\}$, where $r_1 = r_p\{m_p, n_p\}$, $r_2 = r_q\{m_q, n_q\}$, $n > 1$, $\mathcal{L}(r_1) \cap \mathcal{L}(r_2) \neq \emptyset$	$\delta_8 = \text{\@}([\backslash\text{w}\backslash -] + \backslash\text{.}[\backslash\text{w}\backslash - :] +) + [: /]$	$(QOA_3, ([\backslash\text{w}\backslash -] + \backslash\text{.}[\backslash\text{w}\backslash - :] +) +, [[\backslash\text{w}\backslash -] +, [\backslash\text{w}\backslash - :] +])$
#4	QOA_4	$r = (r_1...r_2r_3)\{m, n\}$, where $r_3 = r_t\{m_t, n_t\}$ is nullable, $r_1 = r_p\{m_p, n_p\}$, $r_2 = r_q\{m_q, n_q\}$, $n > 1$, $\mathcal{L}(r_1) \cap \mathcal{L}(r_2) \neq \emptyset$	$\delta_9 = \text{^}([\text{ab}] + \text{d}[\text{ac}] + \text{e}?) + \text{\$}$	$(QOA_4, ([\text{ab}] + \text{d}[\text{ac}] + \text{e}?) +, [[\text{ab}] +, [\text{ac}] +])$
#5	QOA_5	$r = (r_1r_2...r_3)\{m, n\}$, where $r_1 = r_t\{m_t, n_t\}$ is nullable, $r_2 = r_p\{m_p, n_p\}$, $r_3 = r_q\{m_q, n_q\}$, $n > 1$, $\mathcal{L}(r_2) \cap \mathcal{L}(r_3) \neq \emptyset$	$\delta_{10} = \text{^} (; ([\backslash\text{t}]^* [0-9\text{a-zA-Z}] + = [\backslash\text{x}21-\backslash\text{x}7\text{E}]^*)^*, [\backslash\text{x}21-\backslash\text{x}7\text{E}]^*)^* [\backslash\text{t}]^*) ? \text{\$}$	$(QOA_5, ([\backslash\text{t}]^* [0-9\text{a-zA-Z}] + = [\backslash\text{x}21-\backslash\text{x}7\text{E}]^*)^*, [[0-9\text{a-zA-Z}] +, [\backslash\text{x}21-\backslash\text{x}7\text{E}]^*])$

- Quantified Overlapping Adjacent (QOA)

- The QOA pattern is a quantified regex containing two adjacent overlapping sub-regexes.
- In order to facilitate fixing the pathological regex, we subdivided QOA pattern into five sub-patterns (i.e., QOA1, QOA2, QOA3, QOA4, and QOA5).

Vulnerable Patterns

Table 4: The Sub-pattern, Vulnerability Description, Example Regex, and Results from AnaSLQ of the Pattern SLQ .

No.	Sub-pattern	Vulnerability Description	Example Regex	Returned Triple
#1	SLQ_1	Starting with $r = r_1$, where $r_1 = r_q\{m_q, n_q\}$, and $n_q > 1$	$\delta_{11} = ([A-Z]^+)([A-Z][a-z])$ (CVE-2021-3820)	$(SLQ_1, [A-Z]^+, [A-Z]^+)$
#2	SLQ_2	Starting with $r = r_1r_2$, where $r_1 = r_p\{m_p, n_p\}$ is nullable, $r_2 = r_q\{m_q, n_q\}$, and $n_q > 1$	$\delta_{12} = \backslash\$?[A-Z]^+$ (PrismJS)	$(SLQ_2, \backslash\$?[A-Z]^+, [A-Z]^+)$
#3	SLQ_3	Starting with $r = r_1r_2$, where $r_1 = r_p\{m_p, n_p\}$ is not nullable, $r_2 = r_q\{m_q, n_q\}$, and $\mathcal{L}(r_1) \cap \mathcal{L}(r_2) \neq \emptyset$	$\delta_{13} = \{([\s\S]^+)\}$ (CVE-2021-3777)	$(SLQ_3, \{([\s\S]^+)\}, \{([\s\S]^+)\})$
#4	SLQ_4	Starting with $r = r_1r_2$, where $r_1 = r_p\{m_p, n_p\}$ is not nullable, $r_2 = r_q\{m_q, n_q\}$, $r_q = r_{q_1}r_{q_2}$, r_{q_1} is not nullable, $r_{q_2} = r_t\{m_t, n_t\}$, and $\mathcal{L}(r_{q_2}) \cap \mathcal{L}(r_1) \neq \emptyset$	$\delta_{14} = [ab](ca)^+d$	$(SLQ_4, [ab](ca)^+, [[ab], a^+])$
#5	SLQ_5	Starting with $r = r_1r_2$, where $r_1 = r_p\{m_p, n_p\}$ is not nullable, $r_2 = r_q\{m_q, n_q\}$, $r_q = r_{q_1}r_{q_2}r_{q_3}$, r_{q_1} and r_{q_3} are not nullable, $r_{q_2} = r_t\{m_t, n_t\}$, and $\mathcal{L}(r_{q_2}) \cap \mathcal{L}(r_1) \neq \emptyset$	$\delta_{15} = [ab](ca\{1,2\}da)^+e$	$(SLQ_5, [ab](ca\{1,2\}da)^+, [[ab], a\{1,2\}])$

- Starting with Large Quantifier (SLQ)
 - The SLQ pattern is a regex starting with a sub-regex with a large quantifier.
 - In order to facilitate fixing the pathological regex, we subdivided SLQ pattern into three sub-patterns (i.e., SLQ1, SLQ2, SLQ3, SLQ4, and SLQ5).

Repair Patterns

Table 5: The Repair Patterns for Nested Quantifiers ($\mathcal{N}Q$).

No.	Repair Pattern
τ_1	$\frac{\mathcal{L}(r) = \mathcal{L}(r_p\{m_p \times m, n_p \times n\})}{r \implies (r_p\{m_p \times m, n_p \times n\})\{m, n\}} \quad (\mathcal{N}Q_1)$
τ_2	$\frac{\mathcal{L}(r) = \mathcal{L}((r_0 r_p \{m_p, n_p\} r_2)\{m, n\})}{r \implies (r_0 r_p \{m_p, n_p\} r_2)\{m, n\}} \quad (\mathcal{N}Q_2)$
τ_3	$\frac{\mathcal{L}(r) = \mathcal{L}((r_0(\dots r_p \{m_p, n_p\} \dots) r_2)\{m, n\})}{r_1 \implies (\dots r_p \{m_p, n_p\} \dots)} \quad (\mathcal{N}Q_3)$

- Nested Quantifiers (NQ)
 - The NQ pattern has a redundant quantifier by nature. So to fix NQ pattern, we can remove the redundant quantifier.

Repair Patterns

Table 6: The Repair Patterns for Quantified Overlapping Disjunction (QOD).

No.	Repair Pattern	No.	Repair Pattern	No.	Repair Pattern
τ_4	$\frac{\mathcal{L}(r) = \mathcal{L}((r_1 \dots r_k)\{m, n\})}{\alpha \Rightarrow (r_1 \dots r_k)} (QOD_1)$	τ_5	$\frac{t > 1}{r_{p_1} \Rightarrow r_{p_1} (?! \Phi(r_{p_2}))} (QOD_1)$	τ_6	$\frac{s > 1}{r_{q_1} \Rightarrow r_{q_1} (?! \Phi(r_{q_2}))} (QOD_{\{1,2\}})$
τ_7	$\frac{}{r_{p_1} \Rightarrow (?! \Phi(\alpha_2)) r_{p_1}} (QOD_1)$	τ_8	$\frac{}{r_{q_1} \Rightarrow (?! \Phi(\alpha_1)) r_{q_1}} (QOD_1)$	τ_9	$\frac{scs(\alpha_1) = true, \Sigma_{\alpha_1} \setminus first(r_{q_1}) \neq \emptyset}{\alpha_1 \Rightarrow \Theta(\Sigma_{\alpha_1} \setminus first(r_{q_1}))} (QOD_1)$
τ_{10}	$\frac{scs(\alpha_2) = true, \Sigma_{\alpha_2} \setminus first(r_{p_1}) \neq \emptyset}{\alpha_2 \Rightarrow \Theta(\Sigma_{\alpha_2} \setminus first(r_{p_1}))} (QOD_1)$	τ_{11}	$\frac{}{r_p \Rightarrow r_p (? < ! \Phi(\alpha_1))} (QOD_2)$	τ_{12}	$\frac{}{r_{q_1} \Rightarrow (?! \Phi(r_p)) r_{q_1}} (QOD_2)$
τ_{13}	$\frac{scs(\alpha_1) = true, \Sigma_{\alpha_1} \setminus first(r_p) \neq \emptyset}{\alpha_1 \Rightarrow \Theta(\Sigma_{\alpha_1} \setminus first(r_p))} (QOD_2)$	τ_{14}	$\frac{r_p = r_u \{m_u, n_u\}, scs(r_u) = true, \Sigma_{r_u} \setminus first(r_{q_1}) \neq \emptyset}{r_p \Rightarrow \Theta(\Sigma_{r_u} \setminus first(r_{q_1}))} (QOD_2)$		

- Quantified Overlapping Disjunction (QOD)

- The QOD pattern has multiple matching paths across the overlapping disjunctions for a string.
- We proposed three strategies, namely, deleting one overlapping disjunction, adding a lookahead constraint to one overlapping disjunction, and modifying one overlapping disjunction by subtracting the first set of the other one.

Repair Patterns

Table 7: The Repair Patterns for Quantified Overlapping Adjacency (QOA).

No.	Repair Pattern	No.	Repair Pattern	No.	Repair Pattern
τ_{15}	$\frac{\mathcal{L}(r_p) = \mathcal{L}(r_q)}{r_1 r_2 \Rightarrow r_p \{m_p + m_q, n_p + n_q\}} (QOA_1)$	τ_{16}	$\frac{}{r_1 \Rightarrow r_1 (? < \Phi(r_q))} (QOA_1)$	τ_{17}	$\frac{}{r_2 \Rightarrow (? \Phi(r_p)) r_2} (QOA_1)$
τ_{18}	$\frac{n \leq 1}{n_p \Rightarrow n_\mu, n_q \Rightarrow n_\mu} (QOA_{\{1,2\}})$	τ_{19}	$\frac{scs(r_p) = true, \Sigma_{r_p} \setminus first(r_q) \neq \emptyset}{r_p \Rightarrow \Theta(\Sigma_{r_p} \setminus first(r_q))} (QOA_{\{1,3\}})$	τ_{20}	$\frac{scs(r_q) = true, \Sigma_{r_q} \setminus first(r_p) \neq \emptyset}{r_q \Rightarrow \Theta(\Sigma_{r_q} \setminus first(r_p))} (QOA_{\{1,3\}})$
τ_{21}	$\frac{scs(r_p) = true, \Sigma_{r_p} \setminus first(r_q) \neq \emptyset, m_p \geq 1}{r_1 \Rightarrow r_p \{m_p - 1, n_p - 1\} \Theta(\Sigma_{r_p} \setminus first(r_q))} (QOA_1)$	τ_{22}	$\frac{scs(r_q) = true, \Sigma_{r_q} \setminus first(r_p) \neq \emptyset, m_q \geq 1}{r_2 \Rightarrow \Theta(\Sigma_{r_q} \setminus first(r_p)) r_q \{m_q - 1, n_q - 1\}} (QOA_1)$	τ_{23}	$\frac{}{r_1 \Rightarrow (? \Phi(r_q)) r_1} (QOA_3)$
τ_{24}	$\frac{}{r_2 \Rightarrow r_2 (? < \Phi(r_p))} (QOA_3)$	τ_{25}	$\frac{scs(r_p) = true, \Sigma_{r_p} \setminus first(r_q) \neq \emptyset, m_p \geq 1}{r_1 \Rightarrow \Theta(\Sigma_{r_p} \setminus first(r_q)) r_p \{m_p - 1, n_p - 1\}} (QOA_3)$	τ_{26}	$\frac{scs(r_q) = true, \Sigma_{r_q} \setminus first(r_p) \neq \emptyset, m_q \geq 1}{r_2 \Rightarrow r_q \{m_q - 1, n_q - 1\} \Theta(\Sigma_{r_q} \setminus first(r_p))} (QOA_3)$
τ_{27}	$\frac{m_t = 0}{r \Rightarrow r_1 r_3 r_1 r_t \{1, n_t\} r_3} (QOA_2)$	τ_{28}	$\frac{m_t = 0}{r_2 r_3 \Rightarrow r_2 r_2 r_t \{1, n_t\}} (QOA_4)$	τ_{29}	$\frac{m_t = 0}{r_1 r_2 \Rightarrow r_2 r_1 \{1, n_t\} r_2} (QOA_5)$

- Quantified Overlapping Adjacent (QOA)
 - The QOA pattern contains the corresponding two overlapping adjacencies.
 - We proposed three repair strategies, that is, merging the overlapping adjacencies, adding a lookahead constraint to one overlapping adjacency, and modifying one overlapping adjacency.

Repair Patterns

Table 8: The Repair Patterns for Starting with Large Quantifier (SLQ).

No.	Repair Pattern	No.	Repair Pattern	No.	Repair Pattern
τ_{30}	$\frac{}{r \Rightarrow \tilde{r}} (SLQ_{\{1,2,3,4,5\}})$	τ_{31}	$\frac{}{n_q \Rightarrow n_\mu} (SLQ_{\{1,2,3,4,5\}})$	τ_{32}	$\frac{m_p = 0}{r \Rightarrow r_2 r_p \{1, n_p\} r_2} (SLQ_2)$
τ_{33}	$\frac{}{r_p \Rightarrow (?!\Phi(r_q))r_p} (SLQ_3)$	τ_{34}	$\frac{}{r_q \Rightarrow (?!\Phi(r_p))r_q} (SLQ_3)$	τ_{35}	$\frac{scs(r_p) = true, \Sigma_{r_p} \setminus first(r_q) \neq \emptyset}{r_p \Rightarrow \Theta(\Sigma_{r_p} \setminus first(r_q))} (SLQ_3)$
τ_{36}	$\frac{scs(r_q) = true, \Sigma_{r_q} \setminus first(r_p) \neq \emptyset}{r_q \Rightarrow \Theta(\Sigma_{r_q} \setminus first(r_p))} (SLQ_3)$	τ_{37}	$\frac{}{r_p \Rightarrow r_p (? < !\Phi(r_t))} (SLQ_{\{4,5\}})$	τ_{48}	$\frac{}{r_t \Rightarrow (?!\Phi(r_p))r_t} (SLQ_{\{4,5\}})$
τ_{39}	$\frac{scs(r_p) = true, \Sigma_{r_p} \setminus first(r_t) \neq \emptyset}{r_p \Rightarrow \Theta(\Sigma_{r_p} \setminus first(r_t))} (SLQ_{\{4,5\}})$	τ_{40}	$\frac{scs(r_t) = true, \Sigma_{r_t} \setminus first(r_p) \neq \emptyset}{r_t \Rightarrow \Theta(\Sigma_{r_t} \setminus first(r_p))} (SLQ_{\{4,5\}})$		

- Starting with Large Quantifier (SLQ)
 - The SLQ pattern contains the sub-regex starting with a large quantifier (for SLQ1 and SLQ2) or the overlapping sub-regexes (for SLQ3, SLQ4 and SLQ5).
 - We proposed four strategies, namely, adding a start-of-line anchor, replacing the large quantifier with a small one, adding a lookahead to one sub-regex, and modifying one sub-regex by subtracting the first set of the other one.

Case Exhibition

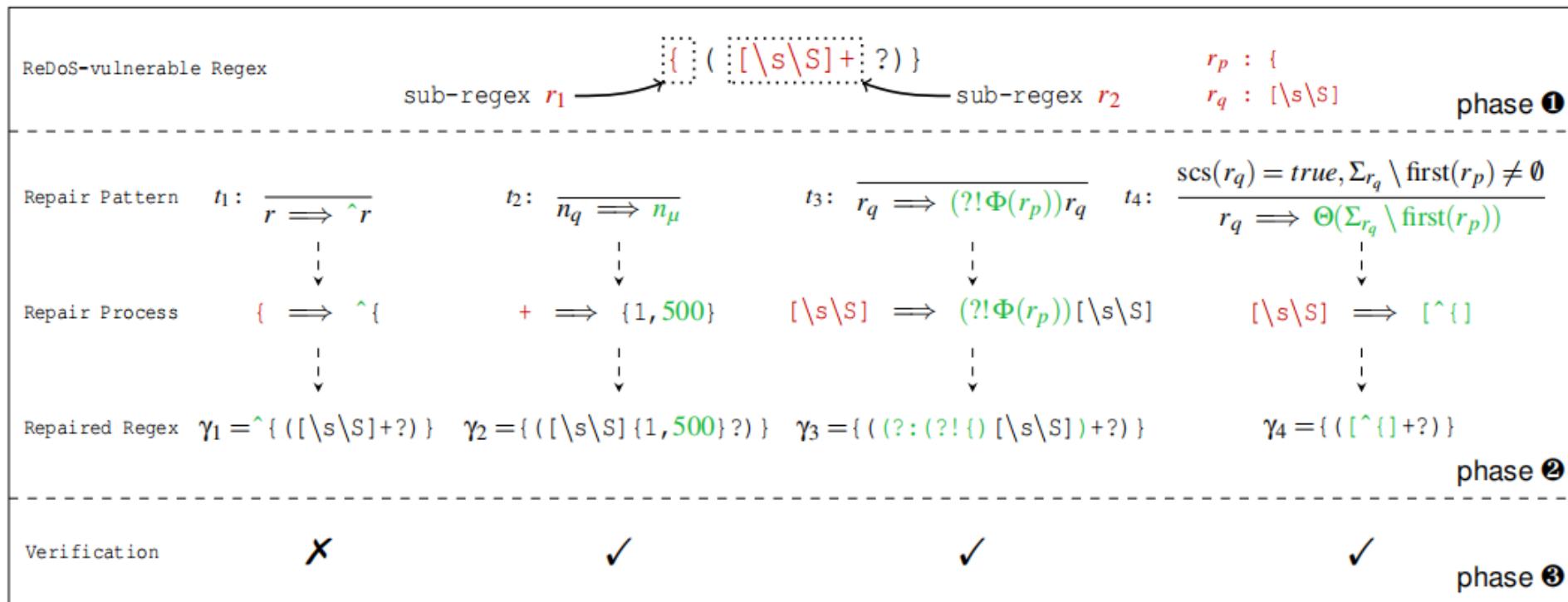


Figure 2: An Example for Repairing the ReDoS-vulnerable Regex $\{ ([\backslash s \backslash S]^+ ?) \}$.

- The NPM package nodejs-tmpl (6,858,130 weekly downloads) used this regex, which aims to match the simple string formatting using `{}`.

Experiment Setup

- Evaluation Datasets
 - Our evaluation was conducted on the ReDoS-vulnerable regexes collected from two widely-used sources: (i) the SOLA-DA benchmark and (ii) real-world CVEs.

Table 9: The ReDoS-vulnerable Regex Sets for Evaluation.

Benchmark	#Regex	Description
SOLA-DA [34]	34	ReDoS-vulnerable regexes in NPM modules found by the Software Lab at TU Darmstadt
CVE [9]	414	ReDoS-vulnerable regexes extracted from 70 ReDoS-related CVEs in recent three years
Total	448	

Experiment Setup

- Evaluation Approaches
 - We selected three state-of-the-art tools belonging to three paradigms, i.e., regex engine substitution (RE2), input length restriction (LLI), and regex repair (FlashRegex).
- Evaluation Metrics
 - A defense is considered successful if it (i) passes all the given test cases, and (ii) is free from ReDoS attack.
 - The success defense rate is calculated by dividing the number of successful defenses by the total number of vulnerable regexes under defense.

Comparing State-of-the-art tools

Table 10: Success Defense Rate Across Automated Tools.

Regex engine substitution

Input length restriction

Input length restriction

Tool	SOLA-DA	CVE	Total
RE2	18 (52.94%)	35 (8.45%)	53 (11.83%)
LLI(100)	22 (64.71%)	45 (10.87%)	67 (14.96%)
LLI(500)	26 (76.47%)	18 (4.35%)	44 (9.82%)
LLI(5000)	26 (76.47%)	19 (4.59%)	45 (10.04%)
FlashRegex	4 (11.76%)	91 (21.98%)	95 (21.20%)
RegexScalpel	33 (97.06%)	410 (99.03%)	443 (98.88%)
#Regex	34	414	448

Our method

- RegexScalpel can effectively defend 98.88% of vulnerable regexes, compared with 21.20% achieved by the best work.

Comparing Maintainers' Repairs

Table 14: Success Defense Rate of the Repairs by Maintainers and RegexScalpel.

Input length restriction

	SOLA-DA	CVE	Total
Manual Repair	14 (66.67%)	305 (77.81%)	319 (77.23%)
RegexScalpel	21 (100%)	388 (98.98%)	409 (99.03%)
#Regex	21	392	413

Our method

- Among the repaired vulnerable regexes handcrafted by the maintainers, only 77.23% are ReDoS free. RegexScalpel outperforms manual fixing, and successfully repairs 99.03% of regexes.

Usefulness to Maintainers

Table 17: Demographics of New Vulnerabilities Repaired by RegexScalpel.

No.	Project	Disclosure Date	CVE ID	#Vuln. Regex
#1	Python	Jan 30th, 2021	CVE-2021-3733	1
#2	NLTK	Sep 5th, 2020	–	1
#3	pylint	Sep 3rd, 2020	–	2
#4	mpmath	Oct 8th, 2021	CVE-2021-29063	1
#5	browserslist	Apr 28th, 2021	CVE-2021-23364	6
#6	code-server	Sep 17th, 2021	CVE-2021-3810	1
#7	ansi-regex	Sep 12th, 2021	CVE-2021-3807	1
#8	nth-check	Sep 17th, 2021	CVE-2021-3803	1
#9	nodejs-tmpl	Sep 15, 2021	CVE-2021-3777	1
#10	jspdf	Feb 12th, 2021	CVE-2021-23353	1
Total				16

- RegexScalpel detected and repaired 16 new ReDoS regexes in ten popular projects.
- All the 16 repairs were accepted by the maintainers and merged into subsequent project releases, resulting in 8 confirmed CVEs.

Semantics Preservation

- We used the following equation to calculate the semantic similarities between the repaired regexes and the original ones.

$$Sim(r_1, r_2) = \frac{|\mathcal{L}(r_1) \cap \mathcal{L}(r_2)|}{|\mathcal{L}(r_1) \cup \mathcal{L}(r_2)|}$$

$$Sim(r_1, r_2) = \lim_{\lambda \rightarrow +\infty} \frac{|\mathcal{L}(r_1)^{\leq \lambda} \cap \mathcal{L}(r_2)^{\leq \lambda}|}{|\mathcal{L}(r_1)^{\leq \lambda} \cup \mathcal{L}(r_2)^{\leq \lambda}|}$$

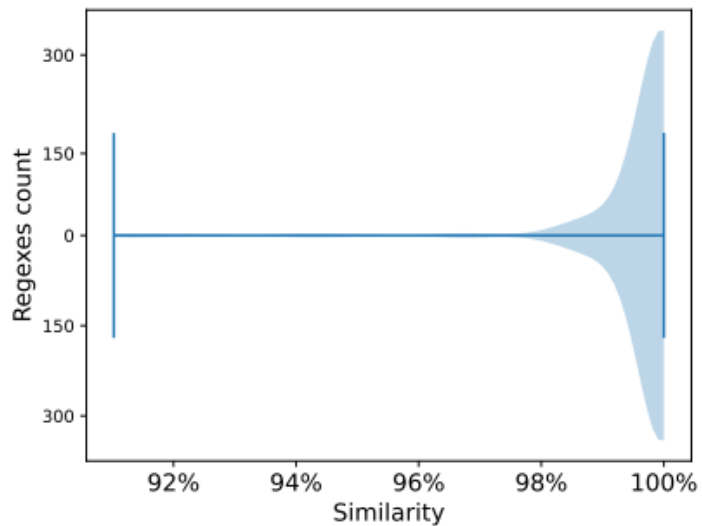


Figure 3: Semantic similarities between the vulnerable regexes and the repaired ones.

- For the benchmarks, most similarities go beyond 98%. On average, the semantic similarity is 99.57%, meaning that the semantics of regexes are well-preserved after the repair.

Summary

- We proposed RegexScalpel, which can defend ReDoS attacks by automatically localizing and repairing vulnerable regexes.
- The evaluation exhibits the remarkable effectiveness of RegexScalpel. It achieves 98.88% successful repair ratio, compared with 21.20% achieved by the best existing work.
- RegexScalpel helped to repair 16 ReDoS vulnerabilities in the ten real-world projects and got confirmed by the maintainers, resulting in 8 confirmed CVEs.
- RegexScalpel can synthesize repaired regexes preserving the semantics of the original ones and keeping the semantics as close as possible to the original ones.



THANKS

Q&A

Presenter: Yecheng Sun

E-mail: sunyc@ios.ac.cn