# From Regular Expressions to Regular Languages

Phúc Phạm

June 10, 2025

# Outline

1. Regex basics

2. Regex tips in practice

3. Regular Languages

4. Introduction to parser

5. Summary

# Regex basics

# What is a regular expression?



What [a-z]{2} [^A-Z]+ (\w*)\?

What is regular expression?

# Regex usecases

## 1. Text search

/found/

## 2. Text validation

example@gmail.com

([a-zA-Z0–9_.+-]+)@[a-zA-Z0–9_.+-]+\.[a-zA-Z0–9_.+-

# Where can I use regex?

4
Unthrifty loveliness why dost thou spend,
Upon thy self thy beauty's legacy?
Nature's bequest gives nothing but doth lend,
And being frank she lends to those are free:
Then beauteous niggard why dost thou abuse,
The bounteous largess given thee to give?
Profitless usurer why dost thou use
So great a sum of sums yet canst not live?
For having traffic with thy self alone,
Thou of thy self thy sweet self dost deceive,
Then how when nature calls thee to be gone,
What acceptable audit canst thou leave?
  Thy unused beauty must be tombed with thee,
  Which used lives th' executor to be.

5
Those hours that with gentle work did frame
The lovely gaze where every eye doth dwell
Will play the tyrants to the very same,
And that unfair which fairly doth excel:
For never-resting time leads summer on
To hideous winter and confounds him there,
Sap checked with frost and lusty leaves quite gone,
Beauty o'er-snowed and bareness every where:

## Find and replace

Find: `\b(up)*on\b`     8 of 380

Replace with:

☐ Match case

☑ Use regular expressions (e.g. \n for newline, \t for tab) Help

☐ Ignore diacritics (e.g. ä = a, E = É, א = א)

Replace    Replace all    Previous    Next

6

# Where can I use regex?
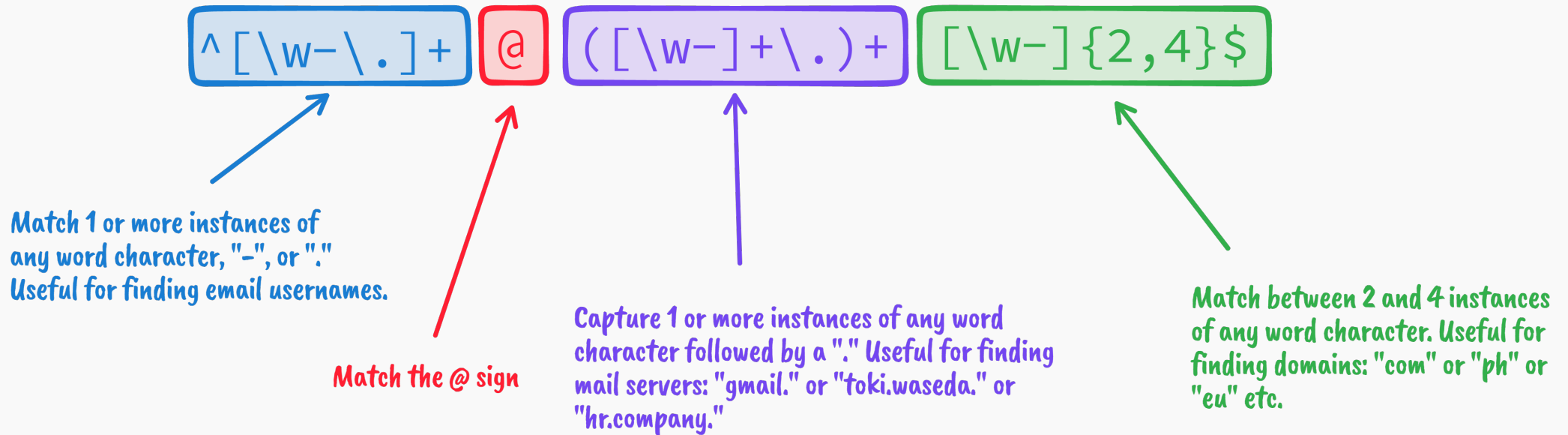
# Where can I use regex?

# Where can I use regex?

```
ppvan at Thinkbook in ~\Downloads
↪rg "\b(up)*on\b" .\t8.shakespeare.small.txt
34:Information on contacting Project Gutenberg to get Etexts, and
48:VERSIONS based on separate sources get new LETTER, shaks10a.txt
167:    (3) You provide or agree to provide on request at no
189:person you received it from.  If you received it on a physical
273:  Thy youth's proud livery so gazed on now,
326:  For never-resting time leads summer on
363:  Attending on his golden pilgrimage:
369:    Unlooked on diest unless thou get a son.
403:    That on himself such murd'rous shame commits.
448:  Borne on the bier with white and bristly beard:
510:  Make war upon this bloody tyrant Time?
513:  Now stand you on the top of happy hours,
623:    Presume not on thy heart when mine is slain,
628:  As an unperfect actor on the stage,
656:  Delights to peep, to gaze therein on thee;
688:  Points on me graciously with fair aspect,
689:  And puts apparel on my tattered loving,
703:  Looking on darkness which the blind do see.
733:  And look upon my self and curse my fate,
739:  Haply I think on thee, and then my state
```

# Regex syntax by example

1

$\texttt{\^{}[\textbackslash w-\textbackslash.]+}$ $\texttt{@}$ $\texttt{([\textbackslash w-]+\textbackslash.)+}$ $\texttt{[\textbackslash w-]\{2,4\}\$}$

Match 1 or more instances of any word character, "-", or "." Useful for finding email usernames.

Match the @ sign

Capture 1 or more instances of any word character followed by a "." Useful for finding mail servers: "gmail." or "toki.waseda." or "hr.company."

Match between 2 and 4 instances of any word character. Useful for finding domains: "com" or "ph" or "eu" etc.

[1]https://regexone.com/

10

# Regex tips in practice

# Use anchor (^$) when possible

Regex: /treature/

treature here, treature there, everywhere treature

---

Regex: /^treature/

treature here, treature there, everywhere treature

---

Regex: /treature$/

treature here, treature there, everywhere treature

# /.*/ is rarely what you want

Regex: /user_id="(.*)"/

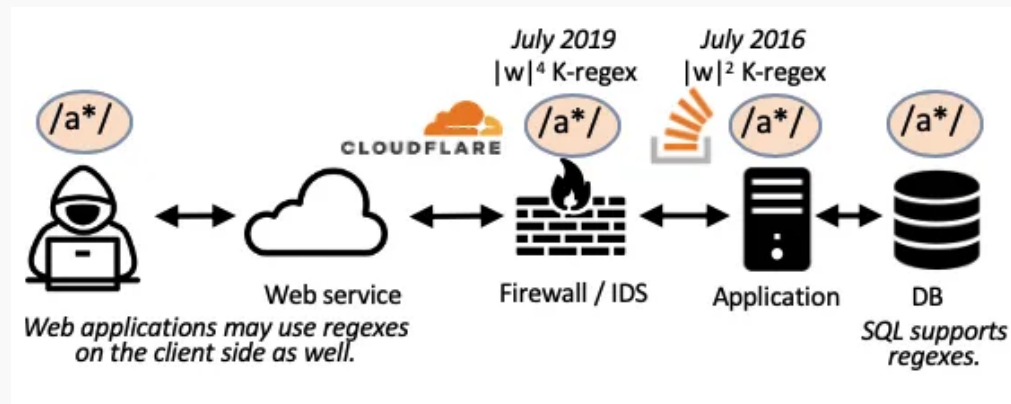2024-01-15 user_id="123" action=login error="Invalid password" ip=192.168.1.100 session="abc123"

---

Regex: /user_id="([^"])"/

2024-01-15 user_id="123" action=login error="Invalid password" ip=192.168.1.100 session="abc123"

# Regex is code

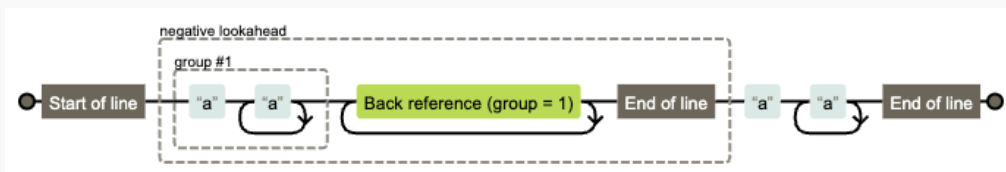- Catastrophic backtracking
- Regular Expression Denial of Service



2

---

# Use visualization tools

regexper.com

regex101.com

# Regular Languages

/a(bb)+a/ $\longrightarrow$

abba
abbbba
abbbbbba
abbbbbbbba
abbbbbbbbbba
.....

$$L = \{a(\mathsf{bb})^n a \mid n \geq 1\}$$

# Finite Automata?

/a(bb)+a/ ←

abba
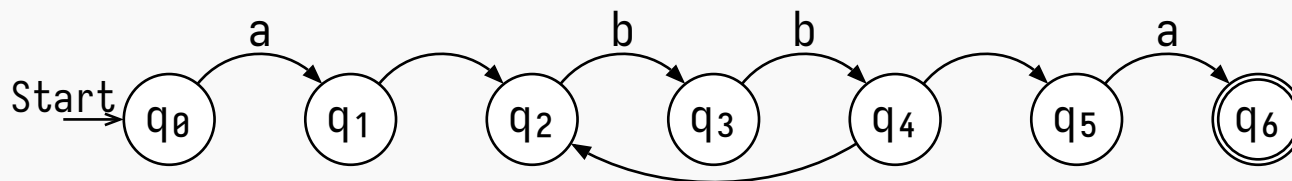abbbba
abbbbbba
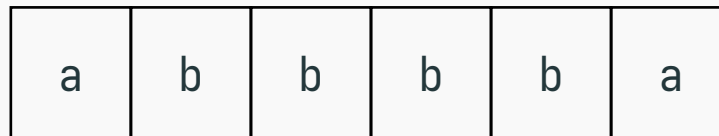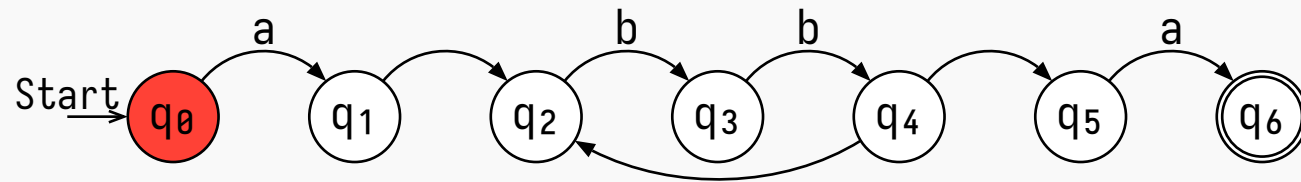abbbbbbbba
abbbbbbbbbba
.....
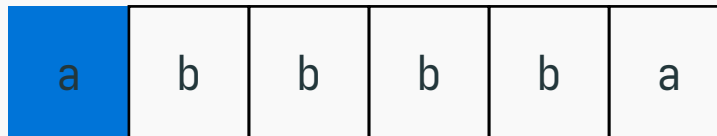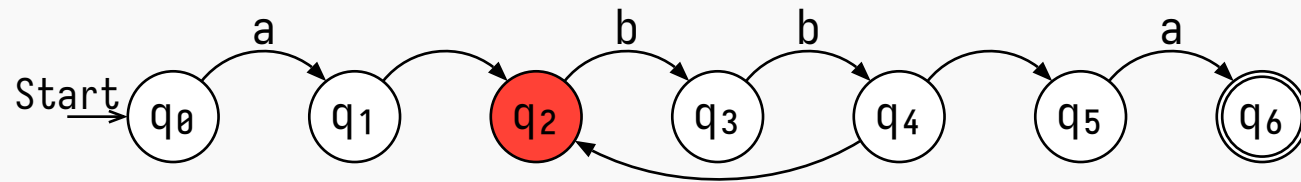
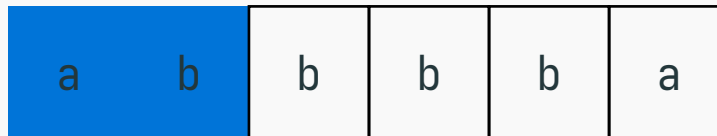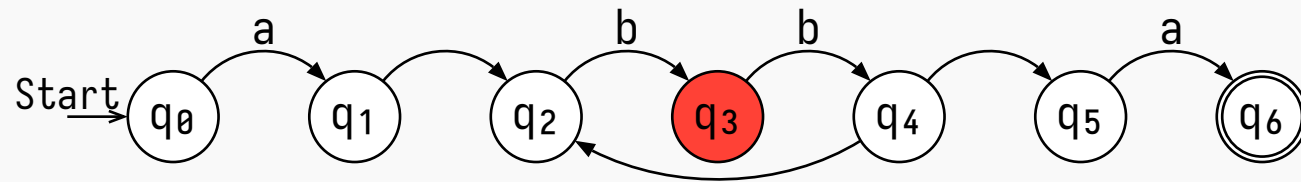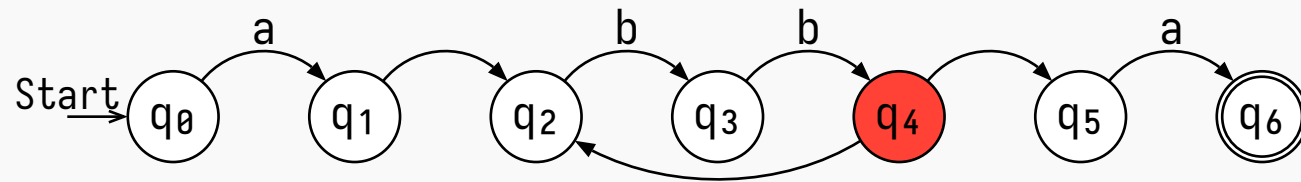

Figure 5: /a(bb)+a/ finite automata

Regex repetitions
(e.g ?*+) causes
backtracking in
regex

$\rightarrow$ O(2^n)
complexity



Figure 6: /a(bb)+a/ state machine

# Catastrophic Backtracking

Pattern: ^(a+)+$

Text: aaaaaaaaaaab

→ ~12.000 steps

REGULAR EXPRESSION                      no match (12,288 steps, 480µs)

/ ^(a+)+$                                                    / gm

TEST STRING

aaaaaaaaaaab

# Introduction to parser

# Why do we need parser?

Regex can't "count"

- html:
- json: {}



Figure 7: NFA can't represent {ab, aabb, aaabbb, ...}

# Regex vs Parsers

Parsers are more powerful and intuitive

### Regex

```
/<p\s+(?:[^>]*\s+)?
class\s*=\s*["\']([^"\']*)["\']
[^>]*>((?:[^<]|<(?!/p>))*)</p>/
```

---

Match all p tag with text

### HTML tree

html
div
image
div
span
p   p

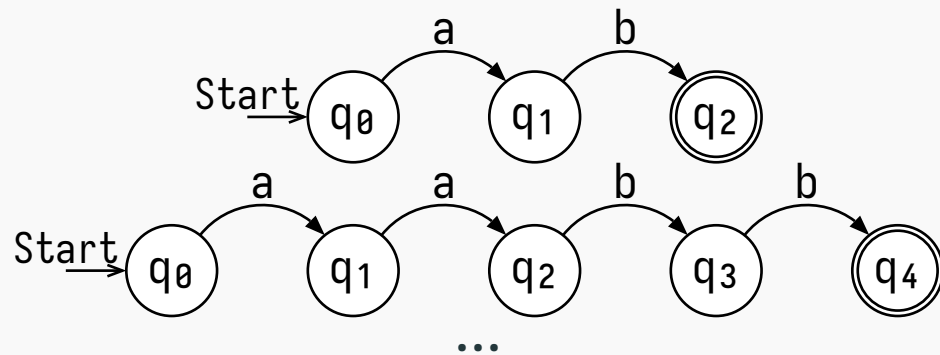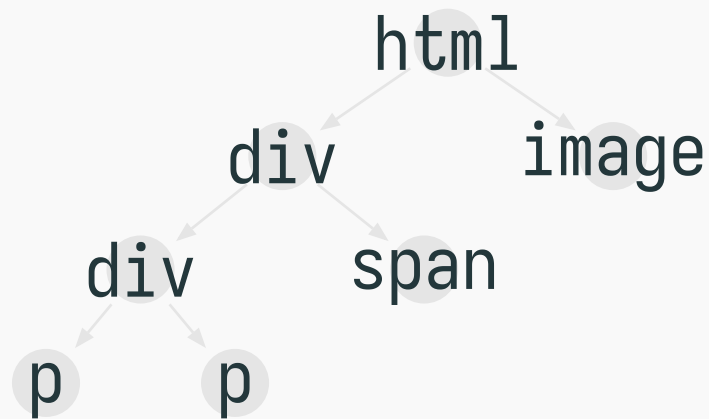# Parsers required more understanding

```python
from lark import Lark
grammar = """
    start: expr
    expr: term ("+" term | "-" term)*
    term: factor ("*" factor | "/" factor)*
    factor: NUMBER | "(" expr ")"
    NUMBER: /[0-9]+/
"""

parser = Lark(grammar)
tree = parser.parse("2 + 3 * 4")
```

Table 1: Parser that match simple math

# Parser usecases

- Compilers & Interpreters
- Data validation: JSON, XML
- Configuration files: YAML, TOML

Recusive Decent

Pratt algorithm

Parser

Parser generator (yacc)

35

# Summary

- Regex is a powerful tools to match "flat" structure extremely fast
- Modern regex is complicated and require careful consideration
- When regex become too hard, consider a parser

Questions?