# Compilers and Formal Languages

Email:          christian.urban at kcl.ac.uk
Slides & Progs:   KEATS

# The Goal of this Module...

## ... you write a compiler

input
program



binary
code

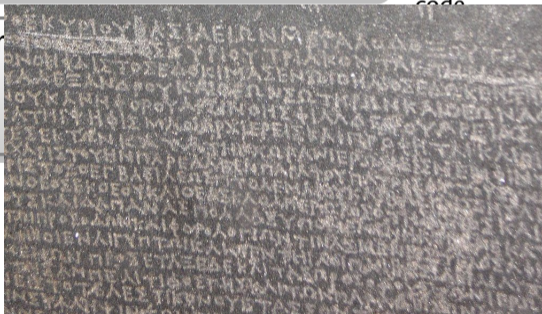lexer → parser → code gen

# The Goal of this Module...

lexer input: a string

```
"read(n);"
```

lexer output: a sequence of tokens

```
key(read) lpar id(n) rpar semi
```

input
program

binary
code



lexer → parser → code gen

# The Goal of this Module...

lexer input: a string

```
"read(n);"
```

lexer output: a sequence of tokens

```
key(read) lpar id(n) rpar semi
```



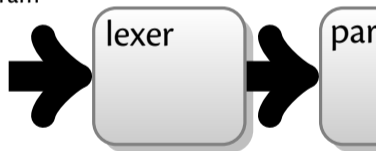lexing ⇒ recognising words (Stone of Rosetta)
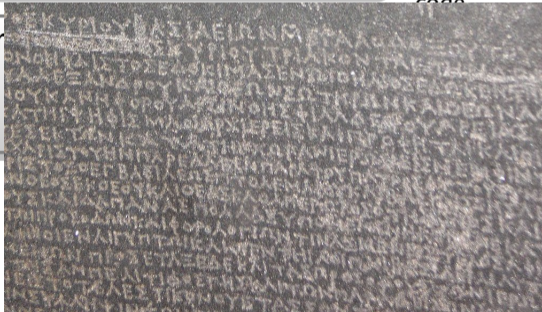
# The Goal of this Module...

lexer input: a string

```
"read(n);"
```

lexer output: a sequence of tokens

```
key(read) lpar id(n) rpar semi
```

input program → lexer → parser → binary code

if ⇒ keyword
iffoo ⇒ identifier



lexing ⇒ recognising words (Stone of Rosetta)

# The Goal of this Module...

parser input: a sequence of tokens

```
key(read) lpar id(n) rpar semi
```

parser output: an abstract syntax tree

```
              read
             / | \
            /  |  \
        lpar   n   rpar
```

inp
pro

binary
code

# The Goal of this Module...

## ... you write a compiler

input
program



binary
code

# The Goal of this Module...

write a compiler

input
pro...

parser

code gen

binary
code

code generation:

```
istore 2
iload 2
ldc 10
isub
ifeq Label2
iload 2
...
```

write a compiler

input
pro

parser

code generation:

```
istore 2
iload 2
ldc 10
isub
ifeq Label2
iload 2
...
```

# The Goal of this Module...

Compiler explorers, e.g.: https://gcc.godbolt.org



input
pro

source ⟶ binary

# The Goal of this Module...

Compiler explorer for Java: `https://javap.yawk.at`

input
pro

```
1   import java.util.*;
2   import lombok.*;
3
4   public class Main {
5       public Main() {
6           int i = 0;
7           i++;
8       }
9   }
```

```
34      Code:
35        stack=1, locals=2, args_size=1
36        start local 0 // Main this
37          0: aload_0
38          1: invokespecial #1
39          4: iconst_0
40          5: istore_1
41        start local 1 // int i
42          6: iinc          1, 1
43          9: return
44        end local 1 // int i
45        end local 0 // Main this
```

## source $\longrightarrow$ byte code

# The Goal of this Module...

## ... you write a compiler

input
program

binary
code

lexer → parser → code gen

# Why Study Compilers?

John Regehr (Univ. Utah, LLVM compiler hacker) 👆

> "...It's effectively a perpetual employment act for solid compiler hackers."

# Why Study Compilers?

John Regehr (Univ. Utah, LLVM compiler hacker) 👈

> **"...It's effectively a perpetual employment act for solid compiler hackers."**

- **Hardware is getting weirder rather than getting clocked faster.**

    "Almost all processors are multicores nowadays and it looks like there is increasing asymmetry in resources across cores. Processors come with vector units, crypto accelerators etc. We have DSPs, GPUs, ARM big.little, and Xeon Phi. This is only scratching the surface."

# Why Study Compilers?

John Regehr (Univ. Utah, LLVM compiler hacker) 👍

> **"...It's effectively a perpetual employment act for solid compiler hackers."**

- **We're getting tired of low-level languages and their associated security disasters.**

  "We want to write new code, to whatever extent possible, in safer, higher-level languages. Compilers are caught right in the middle of these opposing trends: one of their main jobs is to help bridge the large and growing gap between increasingly high-level languages and increasingly wacky platforms."

# Why Bother with Compilers?

**Boeing 777's**: First flight in 1994. They want to achieve triple redundancy for potential hardware faults. 👆

They compile 1 Ada program to

- Intel 80486
- Motorola 68040 (old Macintosh's)
- AMD 29050 (RISC chips used often in laser printers)

using 3 independent compilers.

# Why Bother with Compilers?

**Boeing 777's**: First flight in 1994. They want to achieve triple redundancy for potential hardware faults. 👉

They compile 1 Ada program to



- Intel 80486
- Motorola 68040 (old Macintosh's)
- AMD 29050 (RISC chips used often in laser printers)

using 3 independent compilers.

Airbus uses C and static analysers. Recently started using CompCert.

# What Do Compilers Do?

Remember BF*** from PEP?

| | | |
|---|---|---|
| > | $\Rightarrow$ | move one cell right |
| < | $\Rightarrow$ | move one cell left |
| + | $\Rightarrow$ | increase cell by one |
| - | $\Rightarrow$ | decrease cell by one |
| . | $\Rightarrow$ | print current cell |
| , | $\Rightarrow$ | input current cell |
| [ | $\Rightarrow$ | loop begin |
| ] | $\Rightarrow$ | loop end |
| | $\Rightarrow$ | everything else is a comment |

# A "Compiler" for BF*** to C

```
>  ⇒  ptr++
<  ⇒  ptr--
+  ⇒  (*ptr)++
-  ⇒  (*ptr)--
.  ⇒  putchar(*ptr)
,  ⇒  *ptr = getchar()
[  ⇒  while(*ptr){
]  ⇒  }
   ⇒  ignore everything else
```

```
char field[30000]
char *ptr = &field[15000]
```

# Another "Compiler" for BF to C

$$
\begin{array}{rcl}
\texttt{>...>} & \Rightarrow & \texttt{ptr += n} \\
\texttt{<...<} & \Rightarrow & \texttt{ptr -= n} \\
\texttt{+...+} & \Rightarrow & \texttt{(*ptr) += n} \\
\texttt{-...-} & \Rightarrow & \texttt{(*ptr) -= n} \\
\texttt{.} & \Rightarrow & \texttt{putchar(*ptr)} \\
\texttt{,} & \Rightarrow & \texttt{*ptr = getchar()} \\
\texttt{[} & \Rightarrow & \texttt{while(*ptr)\{} \\
\texttt{]} & \Rightarrow & \texttt{\}} \\
& \Rightarrow & \text{ignore everything else}
\end{array}
$$

```
char field[30000]
char *ptr = &field[15000]
```

# A Brief Compiler History

- Turing Machines, 1936 (a tape as memory)
- Regular Expressions, 1956
- The first compiler for COBOL, 1957 (Grace Hopper)

- But surprisingly research papers are still published nowadays
- "Parsing: The Solved Problem That Isn't" 👉



Grace Hopper

(she made it to David Letterman's Tonight Show 👉)

# Some Housekeeping

**Exams will be online:**

- final exam in January (30%)
- mid-term shortly after Reading Week (10%)

- weekly engagement (10%)

# Some Housekeeping

**Exams will be online:**

- final exam in January (30%)
- mid-term shortly after Reading Week (10%)

- weekly engagement (10%)

**Weekly Homework (optional):**

- uploaded on KEATS, send answers via email, responded individually
- **all** questions in the exam and mid-term will be from the HW!!

# Some Housekeeping

**Coursework (5 accounting for 45%):**

- matcher (5%)
- lexer (8%)
- parser / interpreter (10%)
- JVM compiler (10%)
- LLVM compiler (12%)

# Some Housekeeping

**Coursework (5 accounting for 45%):**

- matcher (5%)
- lexer (8%)
- parser / interpreter (10%)
- JVM compiler (10%)
- LLVM compiler (12%)

  you can use any programming language you like
  (Haskell, Rust)

# Some Housekeeping

**Coursework (5 accounting for 45%):**

- matcher (5%)
- lexer (8%)
- parser / interpreter (10%)
- JVM compiler (10%)
- LLVM compiler (12%)

you can use any programming language you like (Haskell, Rust)

you can use any code I showed you and uploaded to KEATS...**BUT NOTHING ELSE!**

# Some Housekeeping

**Coursework (5 accounting for 45%):**

- matcher (5%)
- lexer (8%)
- parser / interpreter (10%)
- JVM compiler (10%)
- LLVM compiler (12%)

you can use any programming language you like (Haskell, Rust)
you can use any code I showed you and uploaded to KEATS...**BUT NOTHING ELSE!**

# Lectures 1 - 5

transforming strings into structured data

## Lexing      based on regular expressions

(recognising "words")

## Parsing

(recognising "sentences")



Stone of Rosetta

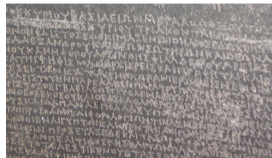# Lectures 1 - 5

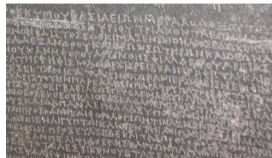transforming strings into structured data

## Lexing

based on regular expressions

(recognising "words")

## Parsing

(recognising "sentences")



Stone of Rosetta

# Lectures 5 - 10

code generation for a small imperative and a small functional language

## Interpreters

(directly runs a program)

## Compilers

(generate JVM code and LLVM-IR code)

# Familiar Regular Expresssions

```
[a-z0-9_\.-]+ @ [a-z0-9\.-]+ . [a-z\.]{2,6}
```

| | |
|---|---|
| `re*` | matches 0 or more times |
| `re+` | matches 1 or more times |
| `re?` | matches 0 or 1 times |
| `re{n}` | matches exactly n number of times |
| `re{n,m}` | matches at least n and at most `m` times |
| `[...]` | matches any single character inside the brackets |
| `[^...]` | matches any single character not inside the brackets |
| `a-z A-Z` | character ranges |
| `\d` | matches digits; equivalent to `[0-9]` |
| `.` | matches every character except newline |
| `(re)` | groups regular expressions and remembers the matched text |

# Some "innocent" examples

Let's try two examples

      `(a*)*b`          `[a?]{n}[a]{n}`

# Some "innocent" examples

Let's try two examples

      `(a*)*b`           `[a?]{n}[a]{n}`

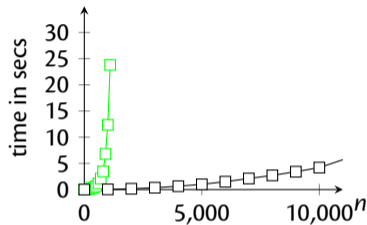and match them with strings of the form

      a, aa, aaa, aaaa, aaaaa, $\underbrace{a...a}_{n}$

# Why Bother with Regexes?



Ruby, Python, Java 8

Us (after next lecture)

`[a?]{n}[a]{n}`:

`(a*)*b`:

matching with strings a...a
          n

# Incidents

- a global outage on 2 July 2019 at **Cloudflare** (first one for six years)

```
(?:(?:\"|'|\]|\}|\\|\d|(?:nan|infinity|true|false|
null|undefined|symbol|math)|\`|\-|\+)+[)]*;?((?:\s
|-|~|!|{}|\|\||\+)*.*(?:.*=.*)))
```



It serves more web traffic than Twitter, Amazon, Apple, Instagram, Bing & Wikipedia combined. 👉

- on 20 July 2016 the **Stack Exchange** webpage went down because of an evil regular expression 👉

# Evil Regular Expressions

- Regular expression Denial of Service (ReDoS)

- Some evil regular expressions:

  - `[a?]{n} [a]{n}`
  - `(a*)* b`
  - `([a-z]+)*`
  - `(a + aa)*`
  - `(a + a?)*`

- sometimes also called catastrophic backtracking

- this is a problem for Network Intrusion Detection systems, Cloudflare, StackExchange, Atom editor

- `https://vimeo.com/112065252`

# (Basic) Regular Expressions

Their inductive definition:

$$
\begin{array}{llll}
r & ::= & \mathbf{0} & \text{nothing} \\
  & | & \mathbf{1} & \text{empty string / "" / []} \\
  & | & c & \text{character} \\
  & | & r_1 + r_2 & \text{alternative / choice} \\
  & | & r_1 \cdot r_2 & \text{sequence} \\
  & | & r^* & \text{star (zero or more)}
\end{array}
$$

**(B**

Their

```scala
abstract class Rexp
case object ZERO extends Rexp
case object ONE extends Rexp
case class CHAR(c: Char) extends Rexp
case class ALT(r1: Rexp, r2: Rexp) extends Rexp
case class SEQ(r1: Rexp, r2: Rexp) extends Rexp
case class STAR(r: Rexp) extends Rexp
```

$$
\begin{array}{llll}
r & ::= & \mathbf{0} & \text{nothing} \\
  & | & \mathbf{1} & \text{empty string / } \texttt{""} \text{ / } [] \\
  & | & c & \text{character} \\
  & | & r_1 + r_2 & \text{alternative / choice} \\
  & | & r_1 \cdot r_2 & \text{sequence} \\
  & | & r^* & \text{star (zero or more)}
\end{array}
$$

# Strings

…are lists of characters. For example `"hello"`

$$[h, e, l, l, o] \text{ or just } hello$$

the empty string: $[]$ or `""`

the concatenation of two strings:

$$s_1 @ s_2$$

*foo @ bar = foobar*
*baz @ [] = baz*

# Languages, Strings

- **Strings** are lists of characters, for example
  $$[], abc \qquad \text{(Pattern match: } c :: s)$$

- A **language** is a set of strings, for example
  $$\{[], hello, foobar, a, abc\}$$

- **Concatenation** of strings and languages
  $$foo \ @ \ bar \ = \ foobar$$
  $$A \ @ \ B \ \stackrel{\text{def}}{=} \ \{s_1 \ @ \ s_2 \ | \ s_1 \in A \land s_2 \in B\}$$

# Languages, Strings

- **Strings** are lists of characters, for example

  $[]$, $abc$ (Pattern match: $c :: s$)

- A **language** is a set of strings, for example

  $\{[], hello, foobar, a, abc\}$

  Let
  $A = \{foo, bar\}$
  $B = \{a, b\}$

  $A @ B = \{fooa, foob, bara, barb\}$

- **Concatenation** of strings and languages

  $foo @ bar = foobar$

  $A @ B \stackrel{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in A \wedge s_2 \in B\}$

# Two Corner Cases

$$A \mathbin{@} \{[]\} = ?$$

# Two Corner Cases

$$A \,@\, \{[]\} = \,?$$

$$A \,@\, \{\} = \,?$$

# The Meaning of a Regex

...all the strings a regular expression can match.

$$
\begin{aligned}
L(\mathbf{0}) &\stackrel{\text{def}}{=} \{\} \\
L(\mathbf{1}) &\stackrel{\text{def}}{=} \{[]\} \\
L(c) &\stackrel{\text{def}}{=} \{[c]\} \\
L(r_1 + r_2) &\stackrel{\text{def}}{=} L(r_1) \cup L(r_2) \\
L(r_1 \cdot r_2) &\stackrel{\text{def}}{=} L(r_1) \,@\, L(r_2) \\
L(r^*) &\stackrel{\text{def}}{=}
\end{aligned}
$$

$L$ is a function from regular expressions to sets of strings (languages):

$L : \text{Rexp} \Rightarrow \text{Set}[\text{String}]$

# The Power Operation

- The **nth Power** of a language:

$$A^0 \stackrel{\text{def}}{=} \{[]\}$$
$$A^{n+1} \stackrel{\text{def}}{=} A @ A^n$$

For example

$$A^4 = A @ A @ A @ A \quad (@ \{[]\})$$
$$A^1 = A \quad\quad\quad\quad\quad (@ \{[]\})$$
$$A^0 = \{[]\}$$

# The Meaning of a Regex

$$L(\mathbf{0}) \overset{\text{def}}{=} \{\}$$

$$L(\mathbf{1}) \overset{\text{def}}{=} \{[]\}$$

$$L(c) \overset{\text{def}}{=} \{[c]\}$$

$$L(r_1 + r_2) \overset{\text{def}}{=} L(r_1) \cup L(r_2)$$

$$L(r_1 \cdot r_2) \overset{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in L(r_1) \wedge s_2 \in L(r_2)\}$$

$$L(r^*) \overset{\text{def}}{=}$$

# The Meaning of a Regex

$$L(\mathbf{0}) \stackrel{\text{def}}{=} \{\}$$

$$L(\mathbf{1}) \stackrel{\text{def}}{=} \{[]\}$$

$$L(c) \stackrel{\text{def}}{=} \{[c]\}$$

$$L(r_1 + r_2) \stackrel{\text{def}}{=} L(r_1) \cup L(r_2)$$

$$L(r_1 \cdot r_2) \stackrel{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in L(r_1) \wedge s_2 \in L(r_2)\}$$

$$L(r^*) \stackrel{\text{def}}{=} \bigcup_{0 \le n} L(r)^n$$

# The Star Operation

- The **Kleene Star** of a language:

$$A\star \overset{\text{def}}{=} \bigcup_{0 \leq n} A^n$$

This expands to

$$A^0 \cup A^1 \cup A^2 \cup A^3 \cup A^4 \cup \ldots$$

or

$$\{[]\} \ \cup\ A \ \cup\ A@A \ \cup\ A@A@A \ \cup\ A@A@A@A \cup \ldots$$

# The Meaning of a Regex

$$L(\mathbf{0}) \stackrel{\text{def}}{=} \{\}$$

$$L(\mathbf{1}) \stackrel{\text{def}}{=} \{[]\}$$

$$L(c) \stackrel{\text{def}}{=} \{[c]\}$$

$$L(r_1 + r_2) \stackrel{\text{def}}{=} L(r_1) \cup L(r_2)$$

$$L(r_1 \cdot r_2) \stackrel{\text{def}}{=} \{s_1 @ s_2 \mid s_1 \in L(r_1) \wedge s_2 \in L(r_2)\}$$

$$L(r^*) \stackrel{\text{def}}{=} (L(r))\star$$

# The Meaning of Matching

A regular expression *r* matches a string *s* provided

$$s \in L(r)$$

…and the point of the next lecture is to decide this problem as fast as possible (unlike Python, Ruby, Java)

# Questions

- Say $A = \{[a], [b], [c], [d]\}$.

  How many strings are in $A^4$ ?

# Questions

- Say $A = \{[a], [b], [c], [d]\}$.

  How many strings are in $A^4$ ?

  What if $A = \{[a], [b], [c], []\}$;
  how many strings are then in $A^4$ ?

**Questions?**