

Parsing with First-Class Derivatives

OOPSLA 2016

Jonathan Brachthäuser
Tillmann Rendel
Klaus Ostermann



Some *text* that

~~~

**function f() {**

~~~

interrupts the code.

~~~

}

~~~

+-----+-----+

| Hello | ASCII |

| World | Tables |

+-----+-----+

while x < 10:

• • y += x

• • x += 1

Some *text* that

~~~

```
function f() {
```

~~~

interrupts the code.

~~~

```
}
```

~~~

+	-----	-----	+
	Hello	ASCII	
	World	Tables	
+	-----	-----	+

```
while x < 10:
```

```
    y += x
```

```
    x += 1
```

Some *text* that

~~~

```
function f() {
```

~~~

interrupts the code.

~~~

```
}
```

~~~

+	-----+	-----+
	Hello	ASCII
	World	Tables
+	-----+	-----+

```
while x < 10:
```

```
    y += x
```

```
    x += 1
```

→ Input to content parsers do not form a **continuous substream!**

The quest for more **modular** and **reusable** parser components.



Agenda

- Parsing with derivatives



Agenda

- Parsing with derivatives
- Parser combinators



Agenda

- Parsing with derivatives
- Parser combinators
- Derivation as a new parser combinator



Agenda

- Parsing with derivatives
- Parser combinators
- Derivation as a new parser combinator
- Examples of gained expressivity



Parsing with Derivatives

$$D_a(a \cdot a \cdot b + b) = a \cdot b$$

$$\mathcal{L}(D_c(r)) = \{w \mid cw \in \mathcal{L}(r)\}$$

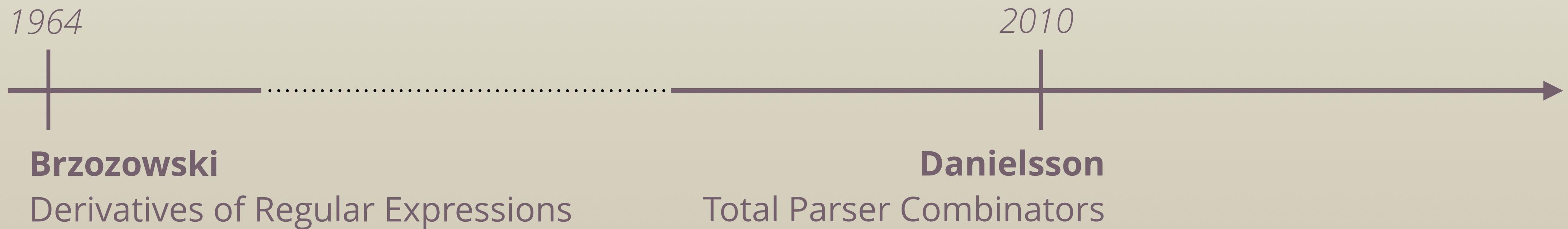
1964



Brzozowski

Derivatives of Regular Expressions

Parsing with Derivatives



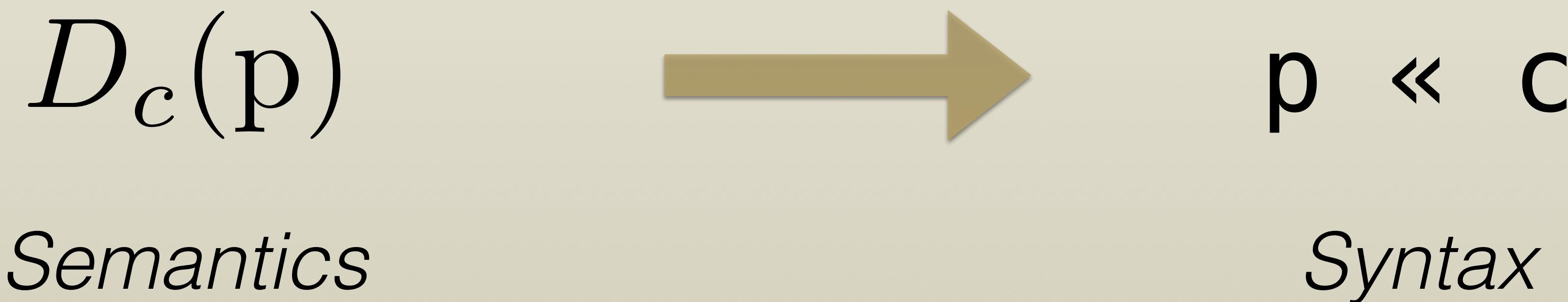
Parsing with Derivatives



Parsing with Derivatives



Parsing with First-Class Derivatives



„Make the (Brzozowski) derivative of a parser available to parser developers.“

Parser Combinators

```
AS ::= ['a' AS]
```

```
def as: Parser[Int] = ...
```



Parser Combinators

```
AS ::= ['a' AS]
```

```
def as: Parser[Int] =  
  ( succeed(0)  
  | ...  
  )
```

Parser Combinators

```
AS ::= ['a' AS]
```

```
def as: Parser[Int] =  
  ( succeed(0)  
  | 'a' ~ as  
  )
```

Parser Combinators

```
AS ::= ['a' AS]
```

```
def as: Parser[Int] =  
  ( succeed(0)  
  | ('a' ~> as) map { n => n + 1 }  
  )
```

Parser Combinators

```
AS ::= ['a' AS]
```

```
def as: Parser[Int] =  
  ( succeed(0)  
  | ('a' ~> as) map { n => n + 1 }  
  )
```

```
> as parse "aaa" → 3
```

Derivation as Parser Combinator

"derive" / "feed"

```
def « [R](p: Parser[R], c: Char) : Parser[R]
```

```
> (as « 'a' « 'a' « 'a') parse "" → 3
```

```
> (as « 'a' « 'a' « 'a') parse "aaa" → 6
```

```
> (as «« "aaa") parse "aaa" → 6
```



Derivation as Parser Combinator

"derive" / "feed"

```
def « [R](p: Parser[R], c: Char) : Parser[R]
```

$$\mathcal{L}(D_c(r)) = \{w \mid cw \in \mathcal{L}(r)\}$$

```
(p « c) parse w = p parse (c + w)
```



Reuse by Selection: An Application

Defined in some other module:

```
val whileStmt: Parser[WhileStmt] = ...
```



Reuse by Selection: An Application

Defined in some other module:

```
val whileStmt: Parser[WhileStmt] = ...
```

Reuse by selection:

```
val untilStmt: Parser[UntilStmt] =  
  "until" ~> (whileStmt << "while") map { ... }
```



Parsing with First-Class Derivatives enables Reuse by Selection



Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

|aaa\n
a\n

"Child-Parser" or "Delegatee"

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

a

<<

a|aa\n
a\n

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aa

<<

a|a\n
a\n

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aaa

<<

aaa|\n
a\n

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aaa

<<

aaa\n
| a\n

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aaaa

<<

aaa\n
a|\n

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aaaa

<<

aaa\n
a\n|

Preprocessing as a Combinator

```
def ignoreNL[R](p: Parser[R]): Parser[R]
```

ignoreNL(as)

as

<<

aaaa

→ 4

<<

aaa\n
a\n|

→ 4

Termination as Parser Combinator

```
def done[R](p: Parser[R]) : Parser[R]
```

```
> done(as << "aaa") parse "" → 3
```

```
> done(as << "aaa") parse "a" → parse error
```

```
> (done(as) << "aaa") parse "" → parse error
```

Preprocessing as a Combinator: Implementation

```
def ignoreNL[R](p: Parser[R]): Parser[R] =  
( done(p)  
| ...  
| ...  
)
```

```
val aLines: Parser[Int] = ignoreNL(as)
```

```
> aLines parse "" → 0
```

Preprocessing as a Combinator: Implementation

```
def ignoreNL[R](p: Parser[R]): Parser[R] =  
  ( done(p)  
  | '\n' ~> ignoreNL(p)  
  | ...  
  )
```

```
> aLines parse "\n" → 0
```

```
> aLines parse "\n\n" → 0
```

Preprocessing as a Combinator: Implementation

```
def ignoreNL[R](p: Parser[R]): Parser[R] =  
  ( done(p)  
  | '\n' ~> ignoreNL(p)  
  | no('\n') flatMap { c => ignoreNL(p << c) }  
  )
```

no: Char => Parser[Char]

Here the input delegation
takes place!



Preprocessing as a Combinator: Implementation

```
def ignoreNL[R](p: Parser[R]): Parser[R] =  
  ( done(p)  
  | '\n' ~> ignoreNL(p)  
  | no('\n') flatMap { c => ignoreNL(p << c) }  
  )
```

```
> aLines parse "a" → 1
```

```
> aLines parse "\na\na\na" → 3
```



Preprocessing as a Combinator: An Application

Unescape derives the delegatee by unescaped characters

```
def unescape[R](p: Parser[R]): Parser[R] =  
  done(p) | "\n" ~> unescape(p << '\n') | ...
```

Preprocessing as a Combinator: An Application

Unescape derives the delegatee by unescaped characters

```
def unescape[R](p: Parser[R]): Parser[R] =  
  done(p) | "\\\n" ~> unescape(p << '\\n') | ...
```

```
val regexString = "" ~> unescape(regex) <~ ""
```



Parsing with First-Class Derivatives enables Local Stream-Preprocessing



Indentation as a Combinator

Simple variant of indentation sensitivity:

```
while • x • < • 10 : ↵  
    • • y • + = • x ↵  
    • • x • + = • 1 ↵
```



Indentation as a Combinator

Simple variant of indentation sensitivity:

```
• • y • + = • x ↘  
• • x • + = • 1 ↘
```



Indentation as a Combinator

Simple variant of indentation sensitivity:

```
indented(aLines)  parse
```

```
  · · aaa ↘  
  · · a ↘
```



Indentation as a Combinator

Simple variant of indentation sensitivity:

indented(aLines) parse

```
  · · aaa ↘  
  · · a ↘
```

=

aLines

parse

```
aaa ↘  
a ↘
```



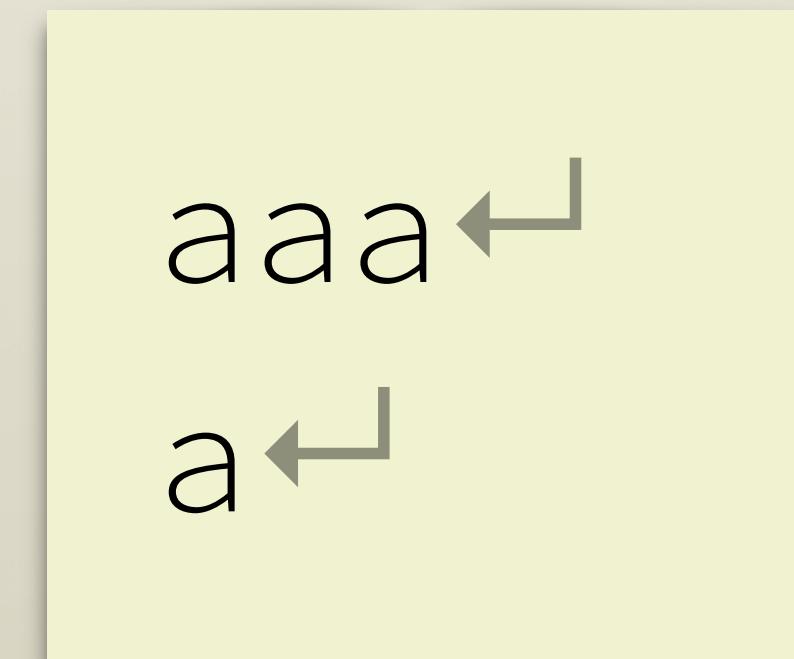
Indentation as a Combinator

Simple variant of indentation sensitivity:

aLines

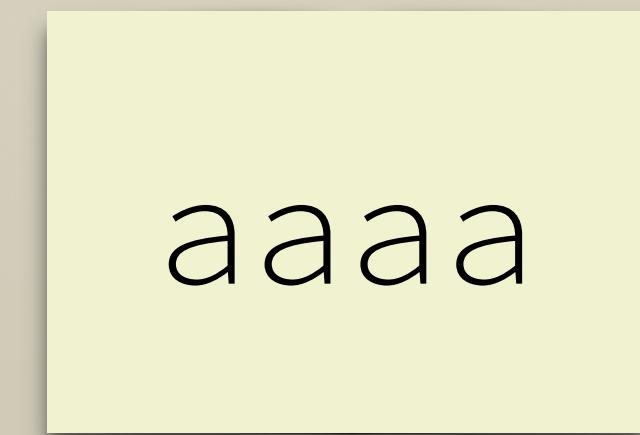
parse

=



as

parse



→ 4



Indentation as a Combinator: Implementation

Simple variant of indentation sensitivity:

```
def indented[T](p: Parser[T]): Parser[T] =  
  done(p) | (space ~ space) ~> readLine(p)
```

```
...aaa\n  
...a\n
```

```
def readLine[T](p: Parser[T]): Parser[T] =  
  ( no('\n') flatMap { c => readLine(p << c) }  
  | '\n'      flatMap { c => indented(p << c) }  
  )
```



Indentation as a Combinator: Implementation

Simple variant of indentation sensitivity:

```
def indented[T](p: Parser[T]): Parser[T] =  
  done(p) | (space ~ space) ~> readLine(p)
```

```
...aaa\n  
...a\n
```

```
def readLine[T](p: Parser[T]): Parser[T] =  
  ( no('\n') flatMap { c => readLine(p << c) }  
  | '\n'      flatMap { c => indented(p << c) }  
  )
```



Indentation as a Combinator: Implementation

Simple variant of indentation sensitivity:

```
def indented[T](p: Parser[T]): Parser[T] =  
  done(p) | (space ~ space) ~> readLine(p)
```

```
...aaa\n  
...a\n
```

```
def readLine[T](p: Parser[T]): Parser[T] =  
  ( no('\n') flatMap { c => readLine(p << c) }  
  | '\n'      flatMap { c => indented(p << c) }  
  )
```

```
> indented(aLines) parse "  a\n  a\n" → 2
```



Indentation as a Combinator: Implementation

Simple variant of indentation sensitivity:

```
def indented[T](p: Parser[T]): Parser[T] =  
  done(p) | (space ~ space) ~> readLine(p)
```

```
• • aaa\n  
• • a\n
```

```
def readLine[T](p: Parser[T]): Parser[T] =  
  ( no('\n') flatMap { c => readLine(p << c) }  
  | '\n'      flatMap { c => indented(p << c) }  
  )
```

```
> indented(indented(aLines))  
parse "      a\n      a\n"
```

→ 2



Parsing with First-Class Derivatives enables Suspending Delegation



In the Paper



- we abstract over recursive delegation patterns
- we implement combinators for the three motivating examples (and more)
- we extend `Indented` to also account for implicit and explicit line joining
- we describe the implementation techniques of our library

Limitations & Future Work

- Improve performance for practicality



Limitations & Future Work

- Improve performance for practicality
- Determine the language class of CFG extended with FCD



Limitations & Future Work

- Improve performance for practicality
- Determine the language class of CFG extended with FCD
- Investigate other forms of derivatives



Limitations & Future Work

- Improve performance for practicality
- Determine the language class of CFG extended with FCD
- Investigate other forms of derivatives
- Combine derivatives with other formalisms, such as PEG



Conclusion

- FCD enable the **reuse** of a parser also in a context, where the input stream to the parser is not a substream of the physical input stream.
- With FCD, some layout features can be **modularly** described as parser combinators.
- Future work could address mentioned limitations and explore other kinds of derivatives

Thank You!

<http://github.com/b-studios/fcd>



End Of Slides

Paradigm: Delegating to Parsers

We abstract over the process of delegating all inputs to a parser

```
def delegate[T](p: P[T]): P[P[T]]
```

```
delegate(p) flatMap done = p
delegate(p) << c       = delegate(p << c)
```

Indentation revisited

With delegation we can reimplement the indentation combinator

```
val line = many(no('\n')) ~ '\n'
```

```
def indented[T](p: P[T]): P[T] =
  ( done(p)
  | (space ~ space) ~> (line &> delegate(p)) >> indented
  )
```



Indentation revisited

With delegation we can reimplement the indentation combinator

```
val line = many(no('`\n')) ~ '`\n'
```

```
def indented[T](p: P[T]): P[T] =  
  ( done(p)  
  | (space ~ space) ~> (line &> delegate(p)) >> indented  
  )  
  "Delimited Delegation"
```



Parsing ASCII tables

```
type Layout = List[Int] // list of column sizes
```

...

```
def delCells[T](l: Layout, cells: List[P[T]]): List[P[P[T]]] =  
  l.zip(cells).map {  
    case (n, p) ⇒ delegateN(n, p).map(p ⇒ p << '\n') <~ '|'  
  }
```

+	-----+	-----+
	Hello	ASCII
	World	Tables
+	-----+	-----+



Delegating to two interleaved Parsers

```
val marker = "~~~\n"
```

```
def inText[R, S](code: P[R]): P[S] => P[(R, S)] =  
( done(code) & done(text)  
| marker  
| not(marker) & line &> delegate(text) >>  
)
```

```
def inCode[R, S](text: P[S]): P[R] => P[(R, S)] = code =>  
( marker  
| not(marker) & line &> delegate(code) >> inCode(text)  
)
```

Some text

~~~

function () {

~~~

Some more text

~~~

}

~~~



Delegating to two interleaved Parsers

```
val marker = "~~~\n"
```

```
def inText[R, S](code: P[R]): P[S] ⇒ P[(R, S)] = text ⇒  
( done(code) & done(text)  
| marker                                ~> inCode(text)(code)  
| not(marker) & line &> delegate(text) >> inText(code)  
)
```

```
def inCode[R, S](text: P[S]): P[R] ⇒ P[(R, S)] = code ⇒  
( marker                                ~> inText(code)(text)  
| not(marker) & line &> delegate(code) >> inCode(text)  
)
```

Implementation of our Library

```
trait P[+R] { p =>
    def results : Res[R]
    def derive  : Elem ⇒ P[R]
    ...
}
```



Implementation of our Library

```
trait P[+R] { p =>
    def results : Res[R]
    def derive  : Elem ⇒ P[R]

    def &[S](q: P[S]) = new P[(R, S)] {
        def results = for(r <- p.results, s <- q.results) yield (r, s)
        def derive  = el => (p derive el) & (q derive el)
    }
    ...
}
```

Implementation of our Library

```
trait P[+R] { p =>
    def results : Res[R]
    def derive  : Elem ⇒ P[R]
    ...
    def «(el: Elem) = p derive el
}
```



Syntax of our Parser Combinator Library

empty lang.

fail : $\forall R. P[R]$

empty word

succeed(r) : $\forall R. R \Rightarrow P[R]$

token select.

acceptIf(f) : $(\text{Elem} \Rightarrow \text{Boolean}) \Rightarrow P[\text{Elem}]$

concat.

p ~ q : $\forall R S. (P[R], P[S]) \Rightarrow P[(R, S)]$

union

p | q : $\forall R. (P[R], P[R]) \Rightarrow P[R]$

map

p ^^ f : $\forall R S. (P[R], R \Rightarrow S) \Rightarrow P[S]$

Boolean Operators

intersection

$$p \& q : \forall R\ S.\ (P[R],\ P[S]) \Rightarrow P[(R,\ S)]$$

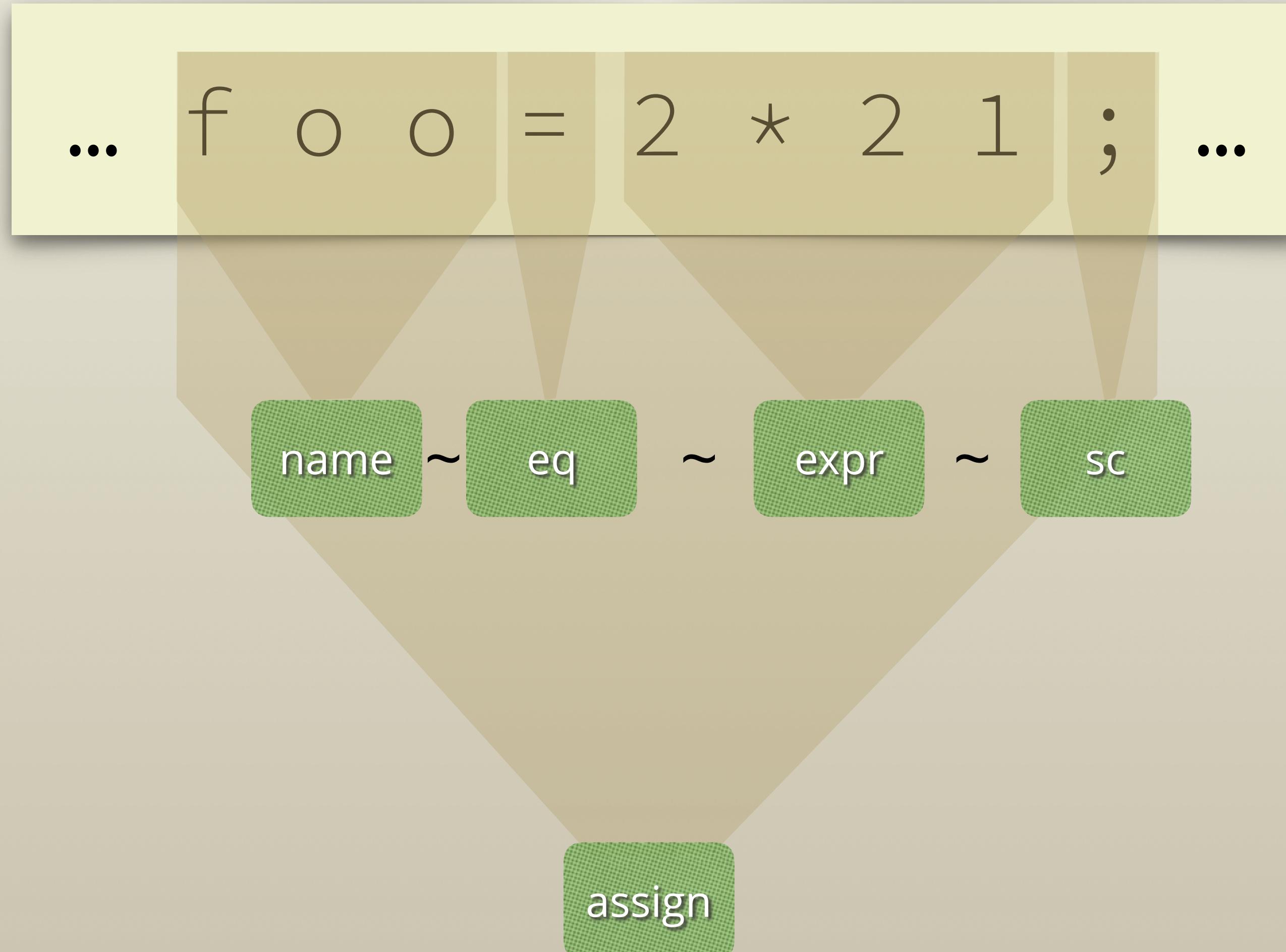
complement

$$\text{not}(p) : \forall R.\ P[R] \Rightarrow P[\text{Unit}]$$

universal

$$\text{always} : P[\text{Unit}]$$

Substream Property



„A parser's virtual input stream corresponds to a continuous substream of the original input stream.“

Limitations implied by the Substream Property

1. Removing from the middle of the stream
2. Extracting interleaved segments of a stream
3. Adding to the stream



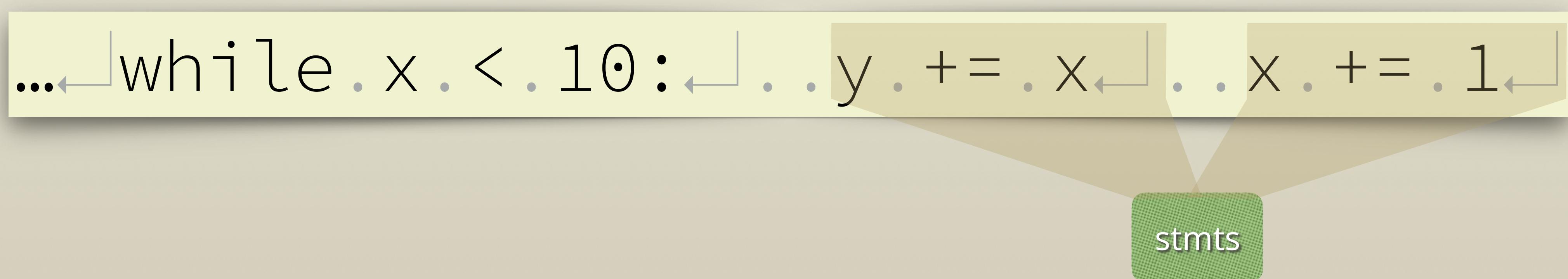
Limitations implied by the Substream Property

1. Removing from the stream



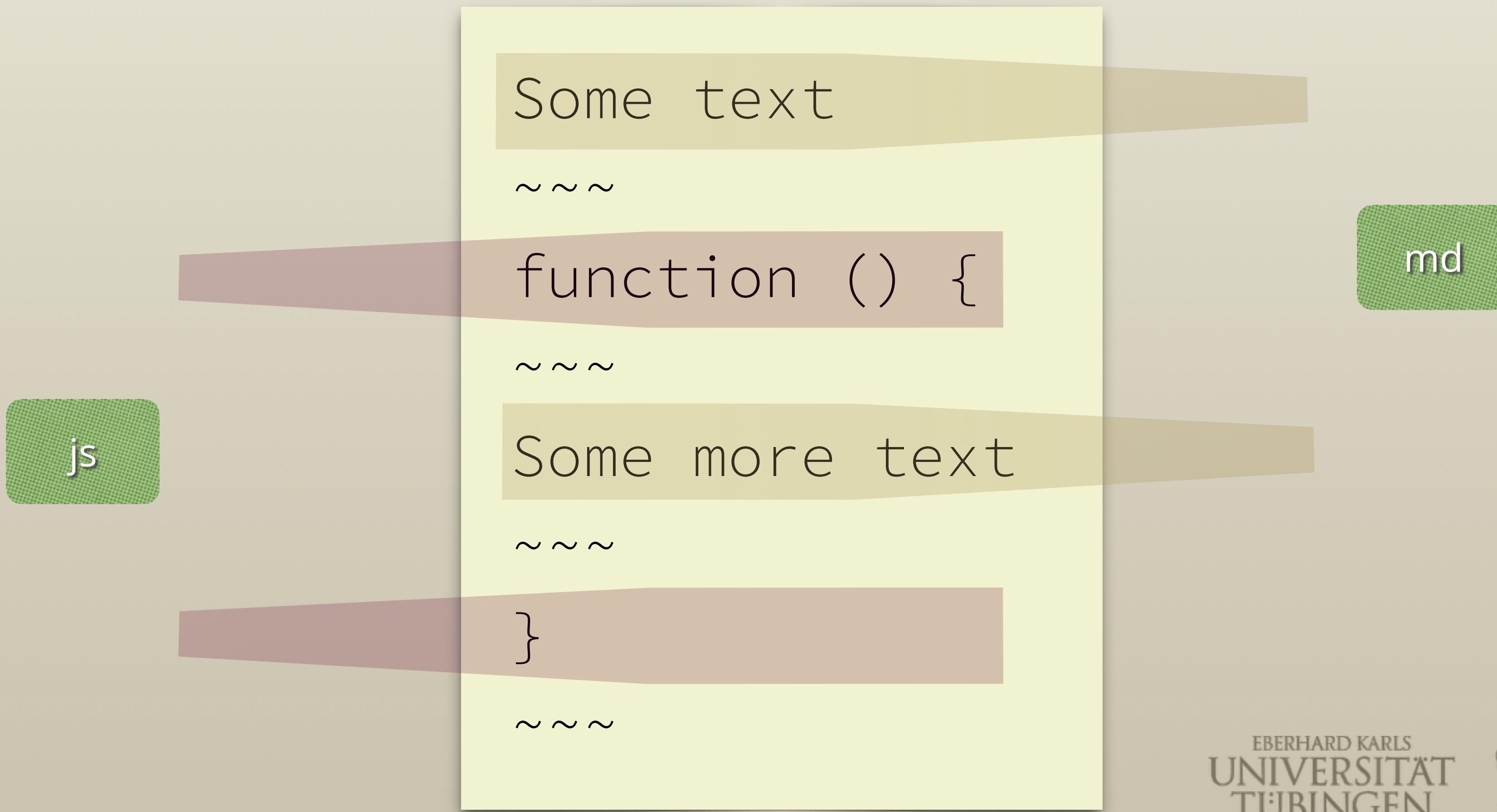
Limitations implied by the Substream Property

1. Removing from the stream



Limitations implied by the Substream Property

2. Extracting interleaved segments of a stream



Limitations implied by the Substream Property

3. Adding to the stream

+-----+	-----+
Hello	ASCII
World	Tables
+-----+	-----+



Limitations implied by the Substream Property

3. Adding to the stream

