

テキスト処理 第5回 (2006-05-23)

田中哲

産業技術総合研究所

情報技術研究部門

`akr@isc.senshu-u.ac.jp`

`http://staff.aist.go.jp/tanaka-akira/textprocess/`

講義の進めかた

- どうも説明不足で速すぎるらしい
- 理解できなければその時点で指摘してくれるとありがたい
- 今日の資料は 40枚くらいあるが、前半部分に焦点をあてて解説する

今日の内容

- mapのレポートの典型的失敗
- egrep.rb -v のレポートの解説
- 簡単な正規表現エンジン
- レポート

mapのレポートの典型的失敗

- 配列の生成
- メソッドの返り値

配列の生成

- `map` は呼ばれるたびに配列を新しく生成しなければならない
- `[...]` は評価のたびに毎回配列を生成する
- `Array.new` も使用できる

メソッドの返り値

- メソッドの返り値はメソッド本体の最後の式の値

```
def m
```

```
  x
```

```
  y
```

```
  z    # z の値が返り値になる
```

```
end
```

- 途中で値を返したいときには return を使える

```
def m
```

```
  return x if y
```

```
  z
```

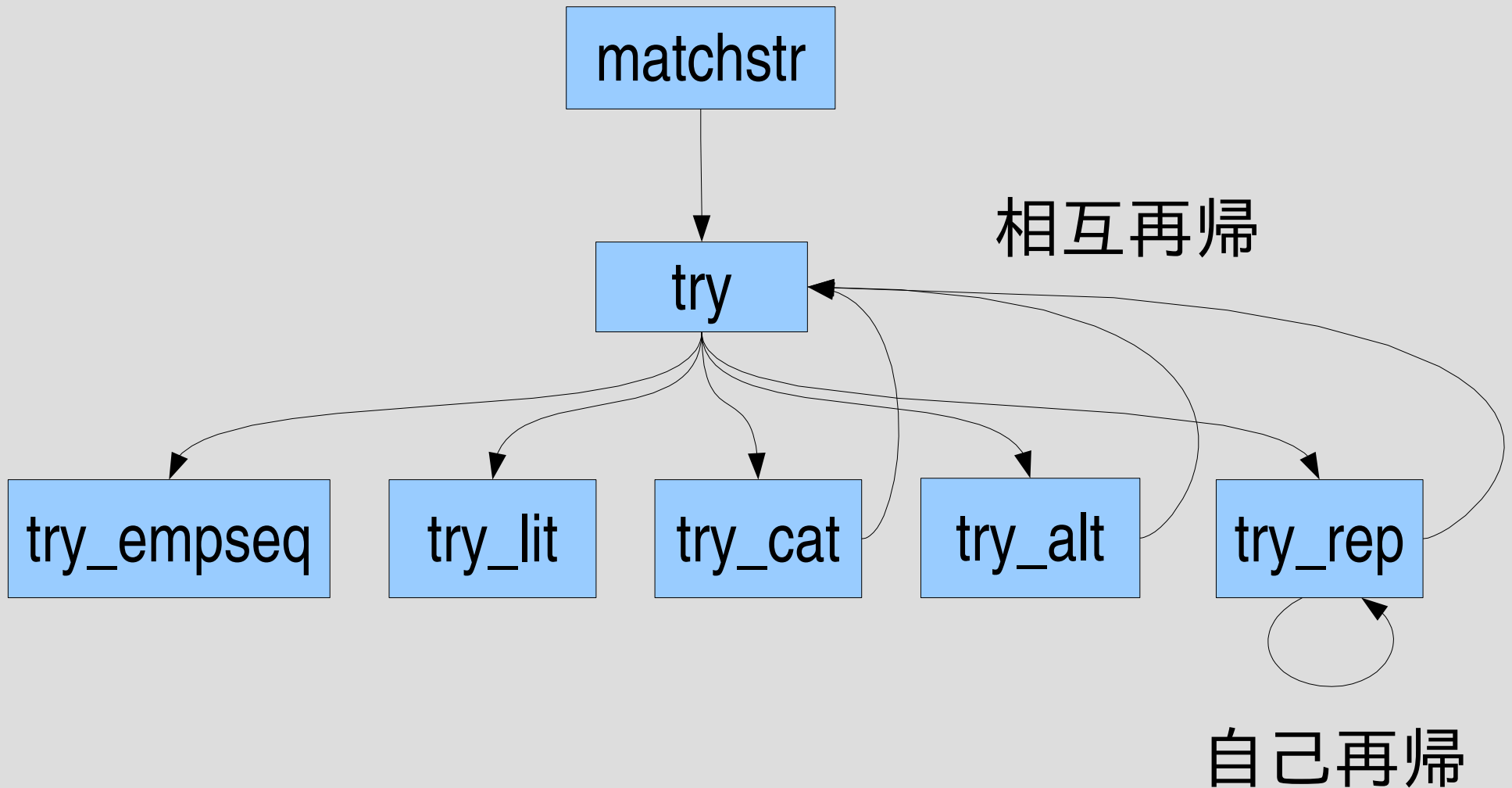
```
end
```

正規表現エンジン

- 6つのメソッドからなる

```
def matchstr(exp, str) ... end # インターフェース
def try(exp, seq, pos) ... end # ディスパッチャ
def try_empseq(seq, pos) ... end # 空文字列 //
def try_lit(sym, seq, pos) ... end # 文字 /x/
def try_cat(e1, e2, seq, pos) ... end # 接続 /XY/
def try_alt(e1, e2, seq, pos) ... end # 選択 /XIY/
def try_rep(exp, seq, pos) ... end # 繰り返し /X*/
```

各メソッドの呼出関係



要素技術

- 文字列の文字への分割 `str.split(//)`
- シンボル
- `case`
- 多重代入
- `&`によるブロックの引きわたり
- 抽象構文木による正規表現
- 再帰

matchstr

```
def matchstr(exp, str)
  result = []
  try(exp, str.split("//"), 0) { |pos|
    result << pos
  }
  result
end
```

try が yield した値を集めて配列として返す
空でない配列になればマッチ

文字列の文字への分割

- String#split(sep)
- sep 区切りで分割して配列になる
- str.split(//) と空文字列で分割すれば文字単位
- 例
 - “”.split(//) #=> []
 - “a”.split(//) #=> [“a”]
 - “abcdef”.split(//) #=> [“a”, “b”, “c”, “d”, “e”, “f”]
- 正規表現エンジンは文字の配列にしてから処理する

matchstrのsplit

```
def matchstr(exp, str)
  result = []
  try(exp, str.split("//"), 0) { |pos|
    result << pos
  }
  result
end
```

tryの仕様

- `try(exp, seq, pos1) { |pos2| ... }`
- `exp` は正規表現の抽象構文木
- `seq` は文字の配列
- `pos1` はマッチを始める位置
- `pos2` はマッチが終わった次の位置
- マッチする可能性すべてについてブロックを呼び出す
 - まったくマッチしなければ呼び出さない
 - 可能性がひとつしかなければ1回だけ呼び出す
 - いろんな可能性があればその数だけ呼び出す

正規表現の抽象構文木

- 正規表現オブジェクトは中身にアクセスできないので違う形で表現する
- 配列、シンボル、文字の組合せ
 - 空文字列 [:empseq] # //
 - 文字 [:lit, "x"] # /x/
 - 接続 [:cat, e1, e2] # /e1e2/
 - 選択 [:alt, e1, e2] # /e1|e2/
 - 繰り返し [:rep, e] # /e*/
- :xxx はシンボル。種類を表現する名前に使用
- 配列以外の表現も考えられる

抽象構文木の例

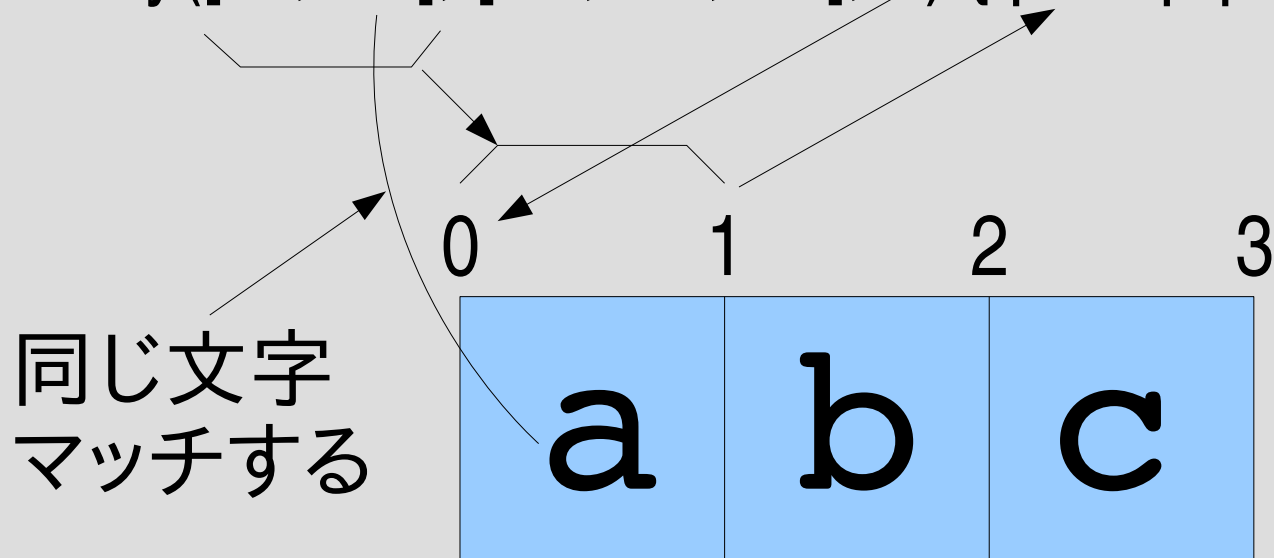
- $[:\text{cat}, [:\text{lit}, \text{"a"}], [:\text{cat}, [:\text{lit}, \text{"b"}], [:\text{lit}, \text{"c"}]]]$
- $[:\text{cat}, [:\text{rep}, [:\text{alt}, [:\text{lit}, \text{"a"}], [:\text{lit}, \text{"b"}]]], [:\text{lit}, \text{"c"}]]]$
- $[:\text{alt}, [:\text{lit}, \text{"a"}], [:\text{empseq}]]]$
- $/\text{abc}/$
- $/(\text{alb})^*\text{c}/$
- $/\text{a}/$

tryの実行例

- `try([:lit, "a"], ["a", "b", "c"], 0) {lposl p pos} # 1`
- `try([:lit, "z"], ["a", "b", "c"], 0) {lposl p pos} # 無し`
- `try([:lit, "a"], ["a", "b", "c"], 1) {lposl p pos} # 無し`
- `try([:rep, [:lit, "a"]], ["a", "b", "c"], 0) {lposl p pos} # 1, 0`
- `try([:rep, [:lit, "a"]], ["a", "b", "c"], 2) {lposl p pos} # 2`
- `try([:rep, [:lit, "a"]], ["a", "b", "c"], 3) {lposl p pos} # 3`
- `try([:rep, [:lit, "a"]], ["a", "a", "a"], 0) {lposl p pos} # 3,2,1,0`

“abc” の 0文字目から /a/

- `try([:lit, “a”], [“a”, “b”, “c”], 0) { |pos| p pos }`

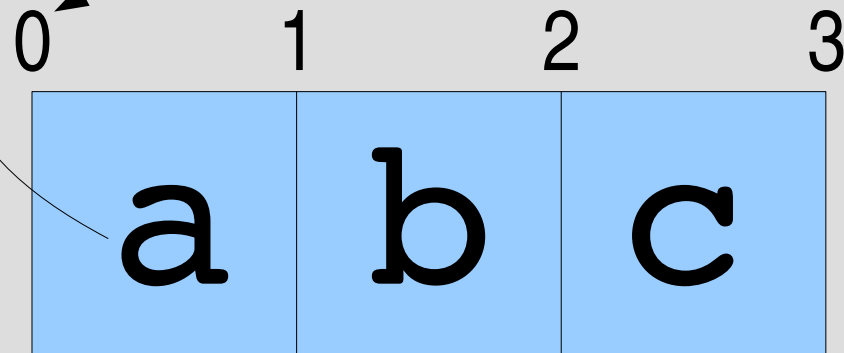


- 1 が一回 yield される

“abc” の 0文字目から /z/

- `try([:lit, “z”], [“a”, “b”, “c”], 0) { |pos| p pos }`

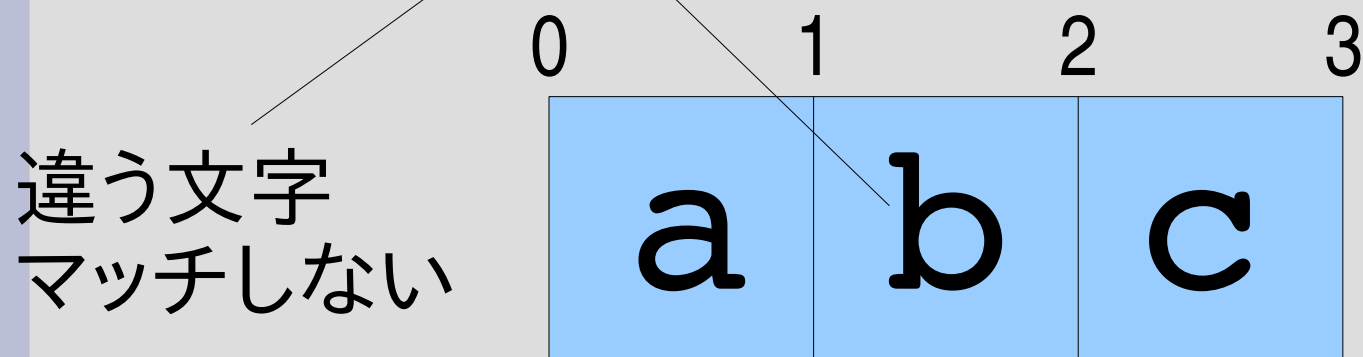
違う文字
マッチしない



- マッチしないので一回も `yield` されない

“abc” の 1文字目から /a/

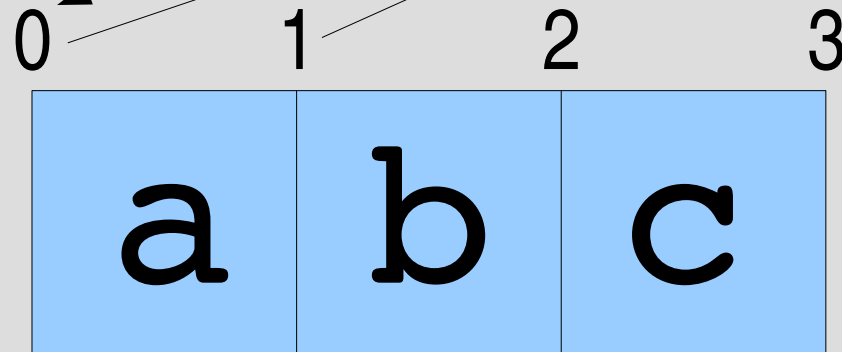
- `try([:lit, “a”], [“a”, “b”, “c”], 1) { |pos| p pos }`



- マッチしないので一回も `yield` されない

“abc” の 0文字目から /a*/

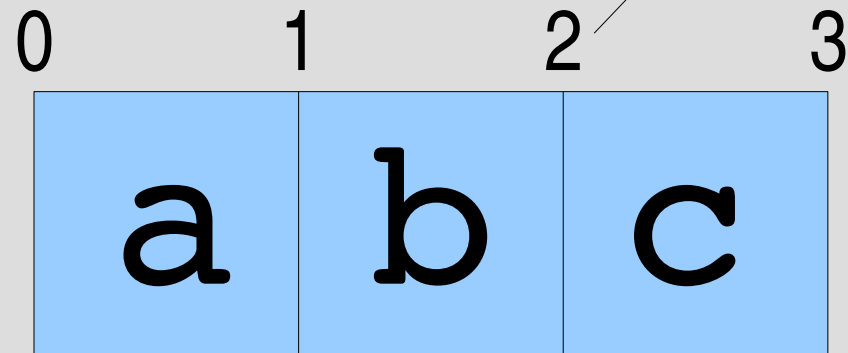
- try([:rep, [:lit, “a”]], [“a”, “b”, “c”], 0) { |pos| p pos }



- /a*/ は “a” と “” にマッチする
- 1 と 0 が順に yield される
- たくさん繰り返した方が先 (最長一致)

“abc” の 2文字目から /a*/

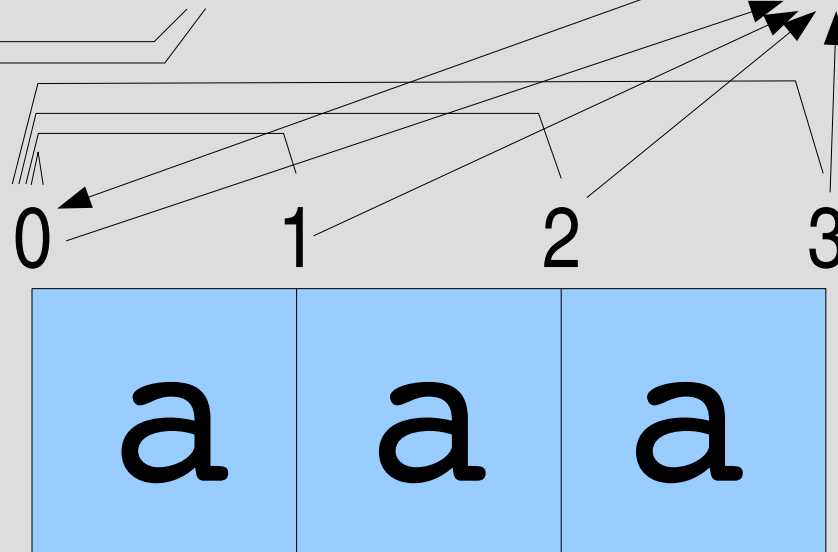
- try([:rep, [:lit, “a”]], [“a”, “b”, “c”], 2) { |pos| p pos }



- /a*/ は “” にマッチする
- 2 が yield される

“aaa” の 0文字目から /a*/

- try([:rep, [:lit, “a”]], [“a”, “a”, “a”], 0) { |pos| p pos }



- /a*/ は “aaa”, “aa”, “a”, “” にマッチする
- 3,2,1,0 の順で yield される

try

```
def try(exp, seq, pos, &b)
  case exp[0]
  when :empseq
    try_empseq(seq, pos, &b)
  when :lit
    _, sym = exp
    try_lit(sym, seq, pos, &b)
  when :cat
    _, e1, e2 = exp
    try_cat(e1, e2, seq, pos, &b)
  when :alt
    _, e1, e2 = exp
    try_alt(e1, e2, seq, pos, &b)
  when :rep
    _, e = exp
    try_rep(e, seq, pos, &b)
  end
end
```

case

- Ruby
 - case 式
when 式
 文
...
else
 文
end
 - 次の選択肢に移ることはない (break 不要)
- C
 - switch (式) {
 case 定数:
 文; break;
...
 default:
 文; break;
}

tryのcase

```
def try(exp, seq, pos, &b)
  case exp[0]
  when :empseq
    try_empseq(seq, pos, &b)
  when :lit
    _, sym = exp
    try_lit(sym, seq, pos, &b)
  when :cat
    _, e1, e2 = exp
    try_cat(e1, e2, seq, pos, &b)
```

```
  when :alt
    _, e1, e2 = exp
    try_alt(e1, e2, seq, pos, &b)
  when :rep
    _, e = exp
    try_rep(e, seq, pos, &b)
  end
end
```

多重代入

- 代入の左辺、右辺には複数の式や配列を書ける

- $a, b = b, a$ #swap



- $a, b, c = \text{array}$

- $a, b, c = [:\text{cat}, e1, e2]$



tryの多重代入

```
def try(exp, seq, pos, &b)
  case exp[0]
  when :empseq
    try_empseq(seq, pos, &b)
  when :lit
    _, sym = exp
    try_lit(sym, seq, pos, &b)
  when :cat
    _, e1, e2 = exp
    try_cat(e1, e2, seq, pos, &b)
```

```
  when :alt
    _, e1, e2 = exp
    try_alt(e1, e2, seq, pos, &b)
  when :rep
    _, e = exp
    try_rep(e, seq, pos, &b)
  end
end
```

&によるブロックの引きわたし

- 例

```
def m1(args, &block) # ブロックを Proc で受け取る
  m2(xxx, &block)    # ブロックを Proc として渡す
end
```
- 以下とだいたい同じ

```
def m1(args)
  m2(xxx) { |v| yield v }
end
```
- 余計なブロックが入らず、短い

tryでのブロックの引きわたり

```
def try(exp, seq, pos, &b)
  case exp[0]
  when :empseq
    try_empseq(seq, pos, &b)
  when :lit
    _, sym = exp
    try_lit(sym, seq, pos, &b)
  when :cat
    _, e1, e2 = exp
    try_cat(e1, e2, seq, pos, &b)
  when :alt
    _, e1, e2 = exp
    try_alt(e1, e2, seq, pos, &b)
  when :rep
    _, e = exp
    try_rep(e, seq, pos, &b)
  end
end
```

tryの内容

- 与えられた exp の種類を case で判別
- exp の内容を多重代入で取り出し
- 種類に応じて try_xxx を呼ぶ

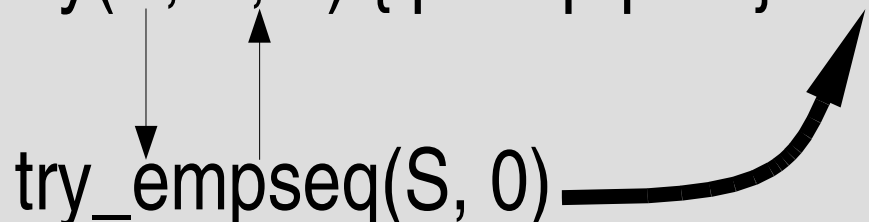
try_empseq

- 空文字列だけ進めて yield
- 空文字列ということはぜんぜん進まないのでものまま yield

```
def try_empseq(seq, pos)
  yield pos
end
```

```
try([:empseq], [], 0) { |pos| p pos } # 0
```

try_empseq の実行

- E=[:empseq] S=[]
 - try(E, S, 0) {lposl p pos } # 0
- try_empseq(S, 0)
- マッチする
- 

try_lit

- 一文字進められれば、進んだ所を yield

```
def try_lit(sym, seq, pos)
  if pos < seq.length && seq[pos] == sym
    yield pos + 1
  end
end
```

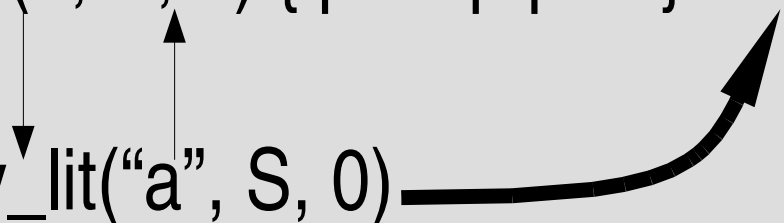
```
try([:lit, "a"], ["a"], 0) { |pos| pos } # 1
```

try_lit の実行

- L=[:lit, "a"] S=["a"]
- try(L, S, 0) { |posl p pos } # 1

try_lit("a", S, 0)

マッチする



try_cat

- e1 を try で進めて、進んだ所から e2 をさらに進める

```
def try_cat(e1, e2, seq, pos, &b)
  try(e1, seq, pos) { |pos2|
    try(e2, seq, pos2, &b)
  }
end
```

```
try([:cat, [:lit, "a"], [:lit, "b"]], ["a", "b"], 0) { |pos1 p pos} # 2
```

try_cat の実行

- L1=[:lit, "a"] L2=[:lit, "b"] C=[:cat, L1, L2]

S=["a", "b"]

- try(C, S, 0) {lpos1 p pos} # 2

try_cat(L1, L2, S, 0)

try_lit("a", S, 0) ← {lpos2 | ...}

マッチする

try_lit("b", S, 1)

マッチする



try_alt

- e1 進めるのを試して、また、e2 進めるのを試す

```
def try_alt(e1, e2, seq, pos, &b)
  try(e1, seq, pos, &b)
  try(e2, seq, pos, &b)
end
```

```
try([:alt, [:lit, "a"], [:lit, "b"]], ["a", "b"], 0) { |pos| p pos } # 1
```

try_alt の実行

- L1=[:lit, "a"] L2=[:lit, "b"] A=[:alt, L1, L2]

S=["a", "b"]

- try(A, S, 0) {lpos| p pos} # 1

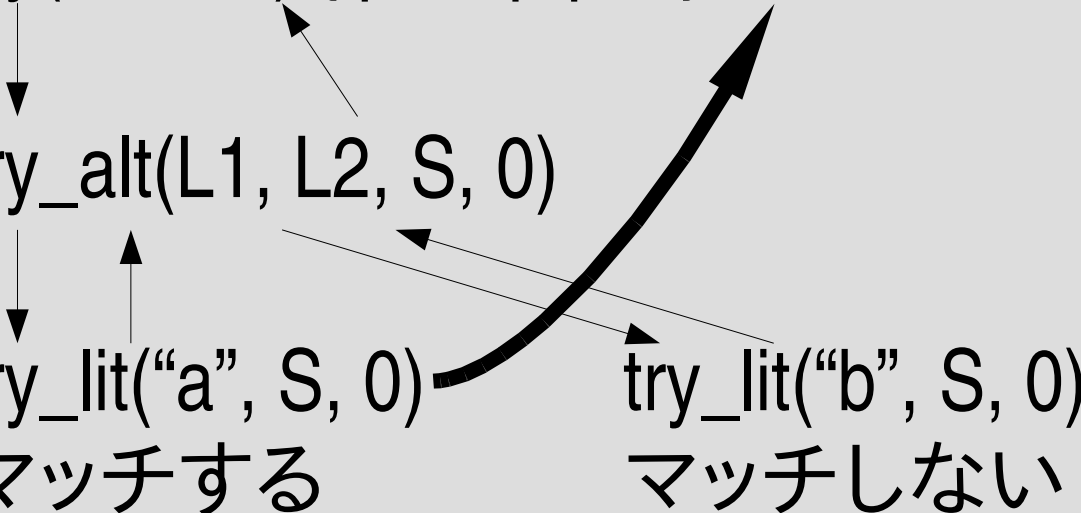
try_alt(L1, L2, S, 0)

try_lit("a", S, 0)

マッチする

try_lit("b", S, 0)

マッチしない



try_rep

- e を進められるだけ進める
 - とりあえず try でひとつ進める
 - ひとつ進めた後に try_rep で進められるだけ進める
- 無限再帰の可能性は気にしない (今は)

```
def try_rep(e, seq, pos, &b)
  try(e, seq, pos) { |pos2|
    try_rep(e, seq, pos2, &b)
  }
  yield pos
end
```

try_rep の実行

- L = [:lit, "a"] R = [:rep, L] S = ["a"]

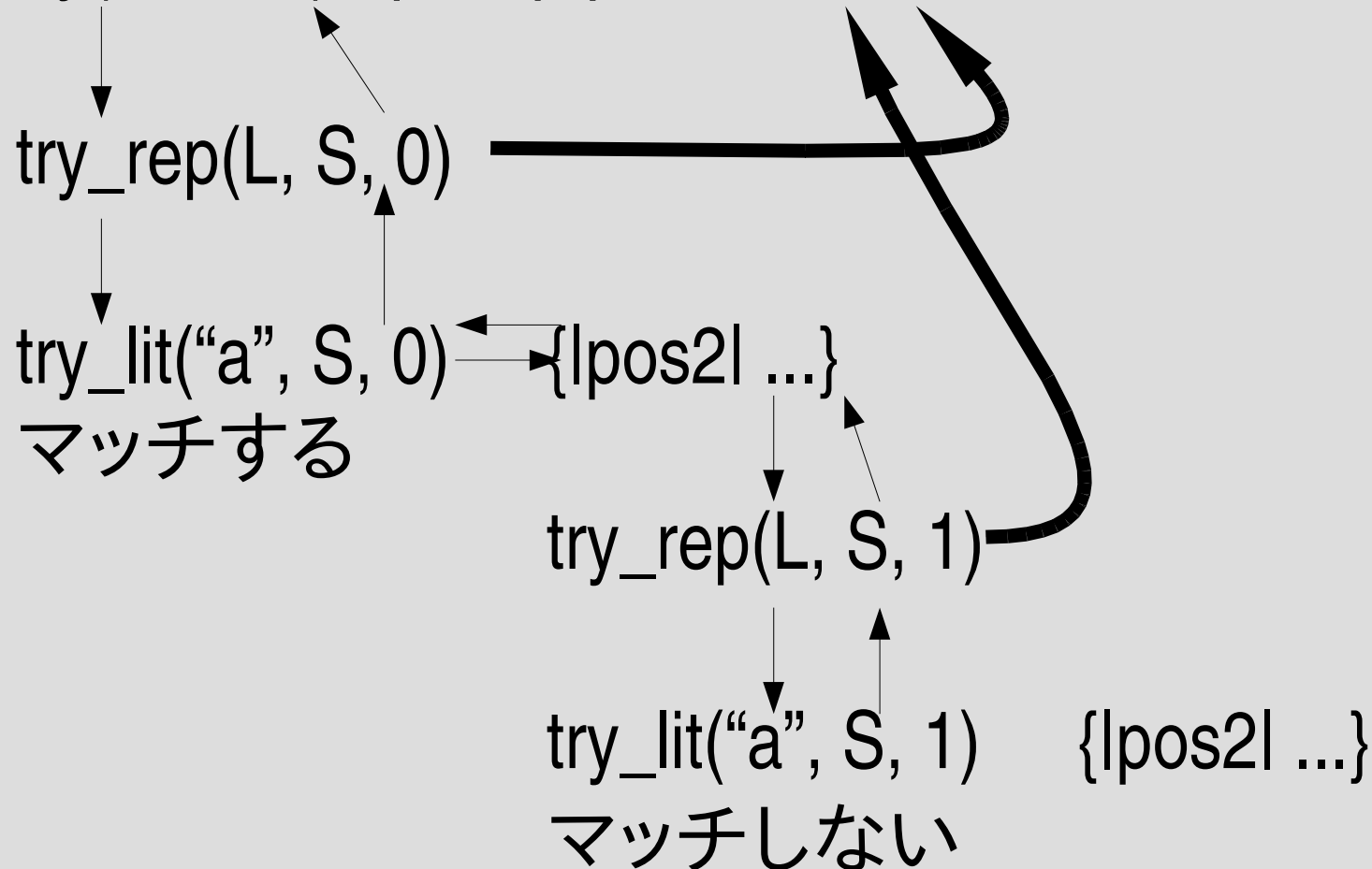
- try(R, S, 0) {lpos1 p pos} # 1, 0

try_rep(L, S, 0)

try_lit("a", S, 0) {lpos2| ...}
マッチする

try_rep(L, S, 1)

try_lit("a", S, 1) {lpos2| ...}
マッチしない



レポート: 正規表現エンジンの使用

- 以下の正規表現を抽象構文木に変換し、今回の正規表現エンジンを用いて動かして動作を解説せよ
 - `/(abc)*abc/`
 - `/abc(abc)*`
 - `/(alb)*abc/`
- ✂切 2006-05-30 16:20
- IT's class
- 拡張子が `txt` なテキストファイルにしてほしい

まとめ

- mapのレポートの典型的失敗
- egrep.rb -v のレポートの解説
- 簡単な正規表現エンジン
 - empseq, lit, cat, alt, rep
- レポートを出した