

A szintaxistól a szemantikáig



Dr. Horpácsi Dániel
ELTE Informatikai Kar
2023-2024-2

A szintaxis és statikus szemantika leírása

- A programozási nyelvek olyan mesterséges nyelvek, amelyeket gépek, számítógépek vezérlésére terveztek
- A nyelv szimbólumai, szavai mondatokat (programokat) alkotnak, amelyek közül a szintaxis határozza meg, melyekhez lehet jelentést rendelni
- A szintaxist általában környezetfüggetlen grammatikával adják meg (BNF formában)

```
<expr>      ::= <literal> | <expr> '+' <expr>  
<literal> ::= '0' | '1' | ... | '9'
```

A környezetfüggetlen szintaxis pontosan definiálja:

- Hogyan lehet a nyelv szimbólumaiból, szavaiból értelmes mondatokat alkotni
- A legtöbb esetben ez egy bővebb nyelvet ad meg, mint amit a fordító ténylegesen elfogad (szintaktikusan jól formált mondatoknak nem feltétlenül van egyértelmű jelentése)
- Milyen nyelvi szerkezeteket lehet használni a nyelvben, és hogyan kell ezeket pontosan jelölni, kombinálni

A grammatika és annak szerkezete általában jól mutatja, hogy

- Hogyan lehet egyszerű kifejezéseket, utasításokat konstruálni
- Hogyan lehet ezeket egymással kombinálni bonyolultabb szerkezetek, absztrakciók alkotásához

A $\mathcal{G} = (\mathcal{T}, \mathcal{N}, \mathcal{P}, \mathcal{S})$ formális grammatika környezetfüggetlen, ha

$$\forall p \in \mathcal{P} : \quad p \equiv A \rightarrow \alpha \quad \text{ahol } A \in \mathcal{N} \text{ és } \alpha \in (\mathcal{T} \cup \mathcal{N})^*$$

- \mathcal{T} : terminális szimbólumok halmaza (a nyelv ábécéje)
- \mathcal{N} : nemterminális szimbólumok halmaza (szintaktikus kategóriák, résznyelvek)
- \mathcal{P} : levezetési szabályok
- $\mathcal{S} \in \mathcal{N}$: kezdőszimbólum (ebből vezetünk le teljes mondatokat)

A mondatok reprezentálhatóak a levezetési fájukkal (konkrét szintaxisfa, 'parse tree'). A szintaxisfa előállítására több ismert és hatékony algoritmus is létezik (pl. LL, LR, LALR, *GLR*).

Példa: környezetfüggetlen grammatika

Legyen $\mathcal{G} = (\mathcal{T}, \mathcal{N}, \mathcal{P}, S)$ környezetfüggetlen grammatika:

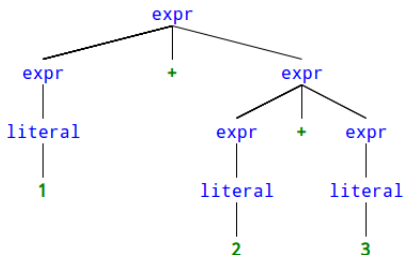
$$\mathcal{T} = \{0, 1, \dots, 9, +\}$$
$$\mathcal{N} = \{expr, literal\}$$
$$S = expr$$
$$\mathcal{P} = \{expr \rightarrow literal, \\ expr \rightarrow expr + expr, \\ literal \rightarrow 0 \\ literal \rightarrow 1 \\ \vdots \\ literal \rightarrow 9\}$$

Ugyanezt a grammatikát BNF-ben jóval olvashatóbban megadhatjuk:

$$\begin{aligned} \langle expr \rangle &::= \langle literal \rangle \mid \langle expr \rangle '+' \langle expr \rangle \\ \langle literal \rangle &::= '0' \mid '1' \mid \dots \mid '9' \end{aligned}$$

```
<expr>      ::= <literal> | <expr> '+' <expr>  
<literal>   ::= '0' | '1' | ... | '9'
```

A grammatika által generált környezetfüggetlen nyelv egy mondata:
1 + 2 + 3



Egyértelmű a fenti grammatika? Egyszerűbben leírható EBNF-fel?

Konkrét szintaxis kontra absztrakt szintaxis

- A **konkrét szintaxis** definiálja, hogyan lehet a tárgyalta programozási nyelven jól formázott mondatokat leírni
- Megadja a nyelvi elemek precíz, konkrét leírását; milyen más elemekből, hogyan konstruálhatunk összetett szerkezeteket
- Szintaktikus elemzőt tudunk ezen definíció alapján készíteni a nyelvhez

`<expr> ::= <expr> '+' <expr>`

- Az **absztrakt szintaxis** csak azt írja le, milyen alapvető szintaktikus kategóriák vannak a nyelvben és azok milyen minták mentén épülnek fel
- Megadja, hogyan épülnek fel az absztrakt szintaxisfák

`expr ::= {add} expr expr`
`add : expr -> expr -> expr`

Az absztrakt szintaxisból kimaradnak például

- Konkrét kulcsszavak
- Precedenciák, asszociativitási szabályok
- Zárójelek
- A szintaktikus elemzést segítő elválasztó és termináló szimbólumok
- Az olvasást (megértést) segítő szimbólumok, jelek

Statikus szemantika

A szintaxisnak nem minden eleme fejezhető ki környezetfüggetlen leírással (néha igen, de rettenetesen bonyolulttá tenné a grammatikát)

- Környezetfüggő szintaxis
- “A szintaxis és a szemantika határa”
- Nem tisztán szerkezeti vagy formai jellemző, de nem is a dinamikus aspektust írja le
- Ez is a szintaxishoz tartozik olyan értelemben, hogy a “mi számít értelmes mondatnak” kérdésre ad választ
- A környezetfüggetlen grammatikáknál kifejezőbb jelölésrendszer kell a leírásához

- Olyan tulajdonságok, amelyeket fordítási időben ellenőrizni lehet
- Tipikus elemei a típusellenőrzés és a névfeloldás

```
1  int i = i;  
2  int i = "foo";  
3  int j = i;  
4  for (int i = 0; i < 10; i++) {  
5      int j = i;  
6  }
```

- (1) Mi az értéke egy deklarálatlan/definiálatlan változónak?
- (2) Mi történik egy változó újradeklarálásakor?
- (2) Mi történik, ha egy `int` változónak szöveget adunk értékül?
- (4) A ciklus fejében lévő `i` változó elrejtí/felüldefiniálja a külső változót?
- (5) A ciklus törzsében lévő `j` változó elrejtí/felüldefiniálja a külső változót?

Amikor feldolgozunk, értelmezünk egy mondatot,

- A szintaktikus elemzés hozza létre a szintaxisfát (majd az AST-t)
- Amelyen a szemantikus elemzés további ellenőrzéseket végez, s kiszűri azokat a mondatokat, amelyeknek biztosan nincs egyértelmű jelentése (pl. olyan programokat, amelyek biztosan futási hibát okoznának)
- A statikus elemzés környezetfüggő: a szintaxisfában egymástól távol eső csúcsokra tesz megkötéseket
- Ehhez szükség van információáramlásra a szintaxisfán belül

- Típusrendszer: levezetési szabályok környezetfüggő típusítéletekre

$$\frac{}{n : \text{Nat}} \quad \frac{e_1 : \text{Nat} \quad e_2 : \text{Nat}}{e_1 + e_2 : \text{Nat}}$$

- Attribútum grammatika: a fa csúcsaira extra információt helyezünk el összetett címkézéssel, az attribútumok értékeire megkötéseket teszünk

```

expr -> num
      { $$ = nat }
| expr '+' expr
  { if($1 == nat && $3 == nat) $$ = nat;
    else error(); }

```

$AG = (\mathcal{G}, \mathcal{A}, \mathcal{R}, \mathcal{C})$ egy attribútum grammatika, ahol

- \mathcal{G} : környezetfüggetlen grammatika, a bázis
- \mathcal{A} : az attribútumok halmaza
- \mathcal{R} : az attribútumszámítási szabályok halmaza ($\mathcal{R}(p)$)
- \mathcal{C} : a feltételek halmaza ($\mathcal{C}(p)$)

Tehát kiterjesztjük a környezetfüggetlen grammatikát attribútumokkal, azok kiszámítási szabályaival, illetve az attribútumokra tett feltételekkel. A levezetett szintaxisfák “dekorálásra” kerülnek, a csúcsokat felcímkezzük az attribútumaikkal.

- Minden terminális és nemterminális szimbólumhoz rendelhetünk attribútumokat
 - $\mathcal{A}(X)$ jelöli az X szimbólum attribútumait
 - $Attr(X)$ (vagy $X.Attr$) jelöli az X szimbólum $Attr$ nevű attribútumát
- A kiszámítási szabályokat a következő formában adjuk meg:

$$Attr1(X) \leftarrow f(Attr2(Y), \dots, AttrN(Z))$$
- Minden attribútumot legfeljebb egy szabály számíthat ki
- A számítás nem hivatkozhat olyan szimbólumok attribútumaira, amelyek nem szerepelnek a szabályban
- (Hasonló megkötések érvényesek a feltételekre)

Az $A.attr$ szintetizált,

- ha van olyan $p \equiv A \rightarrow \alpha$ levezetési szabály,
hogy valamely $r \in \mathcal{R}(p)$ kiszámítja az $A.attr$ attribútumot
- Tehát $A.attr$ olyan szabályban kerül kiszámításra, ahol az A szimbólum a szabály bal oldalán áll
- Az információt felfelé közvetíti a szintaxisfában
- A részfákból számítja ki a fentebbi csúcsok tulajdonságait
- Pl.: névkötések, literálok értéke
- A terminális szimbólumok szintetizált attribútumai speciálisak:
nincsenek részfák, az értékeket egy korábbi elemzési fázis
(lexikális vagy szintaktikus) állítja be

Az $X.attr$ örökölt,

- ha van olyan $p \equiv A \rightarrow \alpha X \beta$ levezetési szabály, ahol $r \in \mathcal{R}(p)$ kiszámítja az $X.attr$ attribútumot
- Azaz X a levezetési szabály jobb oldalán áll, amikor az $X.attr$ kiszámításra kerül
- A szintaxisfában lefelé és keresztben közvetíti az információt
- Pl.: deklarált változók, típusinformációk

Attribútumok nem lehetnek egyszerre szintetizáltak és örököltek is.

Vegyük észre, hogy bármely levezetési fa gyökerében a kezdőszimbólum áll. Következésképp a kezdőszimbólumnak nem lehet örökölt attribútuma (a gyakorlatban persze léteznek megoldások ennek kiküszöbölésére).

A terminális szimbólumok szintetizált attribútumai is speciálisak:

- A terminális szimbólumok a levezetési fa levelei
- Azaz nincsenek részfáik, amikből információt szintetizáljanak
- A gyakorlatban a leveleknek is vannak szintetizált attribútumaik, amelyekben a terminálisok szöveges reprezentációja tárolódik
- Kitüntetett szintetizált attribútum: nem használunk fel a kiszámításakor más attribútumokat
- Enélkül nem lehetne szemantikus elemzést és átírási szemantikát implementálni

Példa: az $a^n b^n c^n$ nyelv grammatikája

$\langle abcSeq \rangle ::= \langle aSeq \rangle \langle bSeq \rangle \langle cSeq \rangle$
 $InSize(\langle bSeq \rangle) \leftarrow Size(\langle aSeq \rangle)$
 $InSize(\langle cSeq \rangle) \leftarrow Size(\langle aSeq \rangle)$

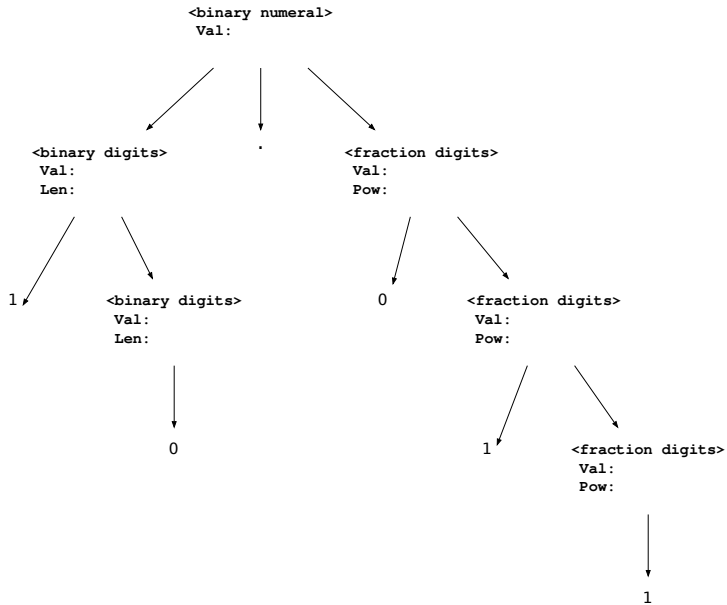
$\langle aSeq \rangle ::= \mathbf{a}$
 $Size(\langle aSeq \rangle) \leftarrow 1$
 $| \quad \langle aSeq \rangle_2 \mathbf{a}$
 $Size(\langle aSeq \rangle) \leftarrow Size(\langle aSeq \rangle_2) + 1$

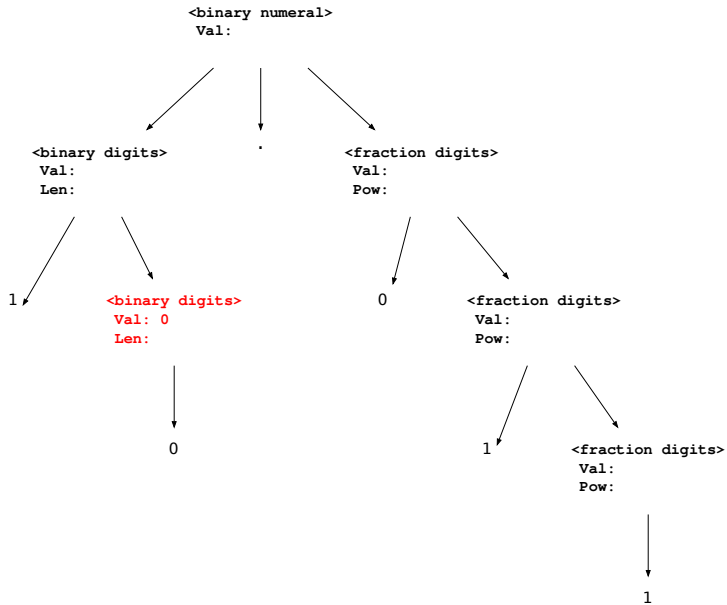
$\langle bSeq \rangle ::= \mathbf{b}$
 $Condition: InSize(\langle bSeq \rangle) = 1$
 $| \quad \langle bSeq \rangle_2 \mathbf{b}$
 $InSize(\langle bSeq \rangle_2) \leftarrow InSize(\langle bSeq \rangle) - 1$

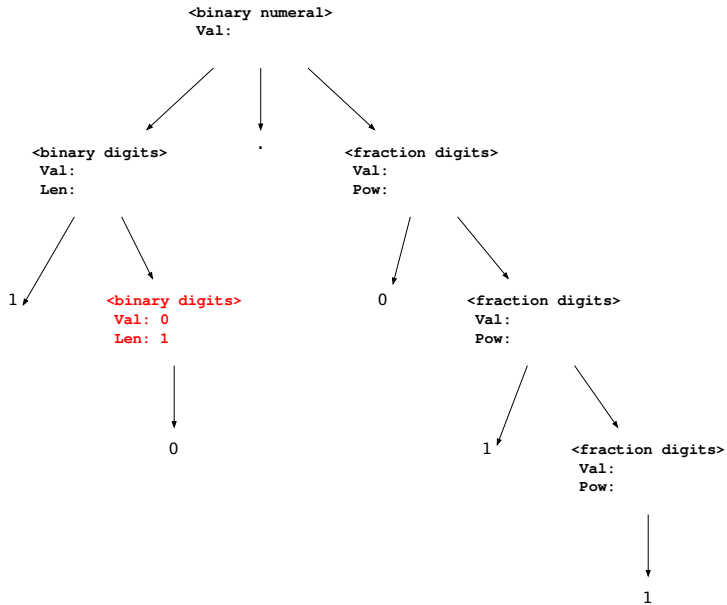
$\langle cSeq \rangle ::= \mathbf{c}$
 $Condition: InSize(\langle cSeq \rangle) = 1$
 $| \quad \langle cSeq \rangle_2 \mathbf{c}$
 $InSize(\langle cSeq \rangle_2) \leftarrow InSize(\langle cSeq \rangle) - 1$

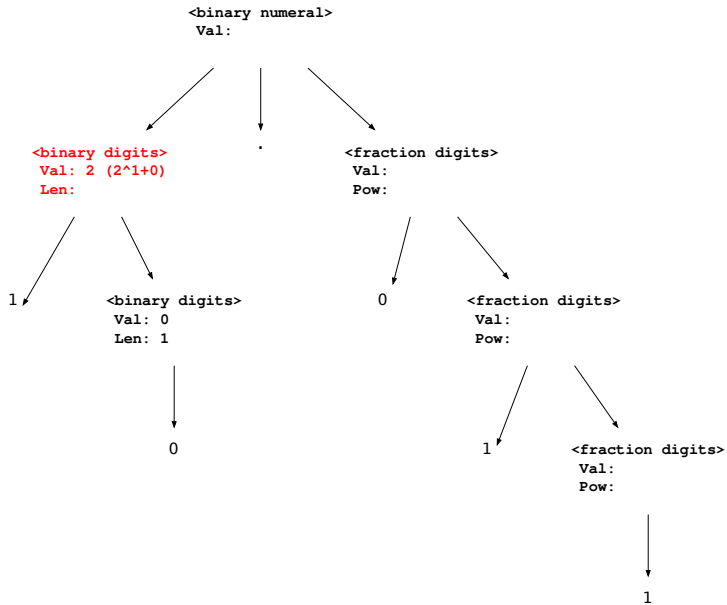
- (Úgy is mondják, hogy a szintaxisfa dekorálása.)
- Az attribútumok, a kiszámítási szabályaik és a feltételek együttesen alkalmasak arra, hogy a nyelv környezetfüggő tulajdonságait leírják
- Amikor egy feltétel kiértékelése hamis értéket eredményez, az szemantikus hibát jelez
- Azonban a kiértékelés nem mindig egyszerű; az attribútumok között függőségek állnak fent, amelyeket figyelembe kell venni
- Annak eldöntése, hogy egy AG szintaxisfája kiértékelhető-e, NP-teljes probléma
- A megfelelő kiszámítási sorrend megtalálására vannak heurisztikák

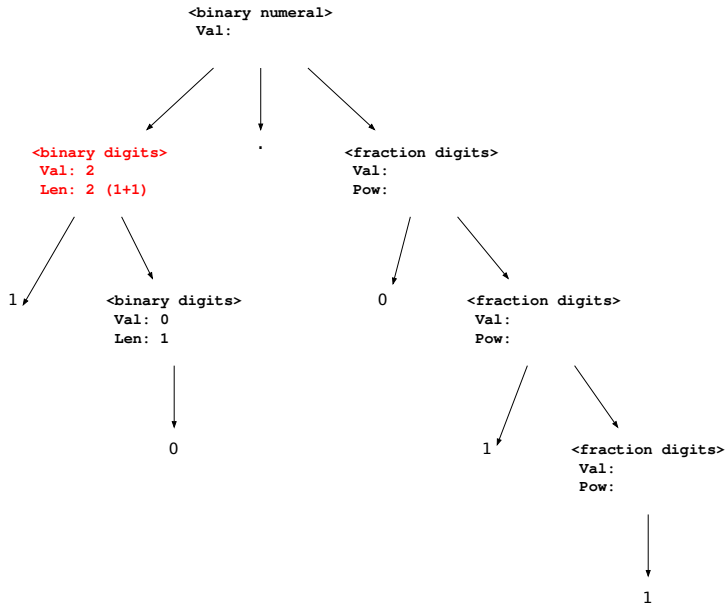
- Tekintsük a *exercises_02.pdf* dokumentumban felírt attribútum grammatikát
- A következő ábrákon látható a “10.011” mondathoz tartozó szintaxisfa attribútumainak egy lehetséges kiértékelése

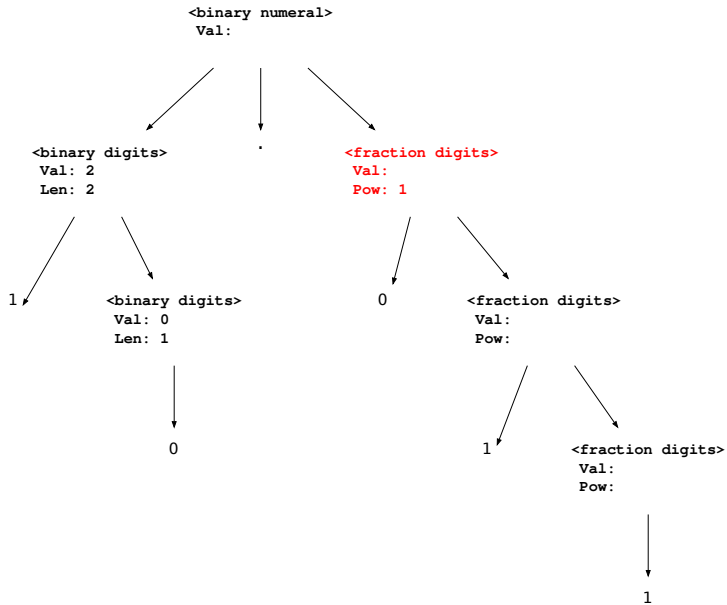


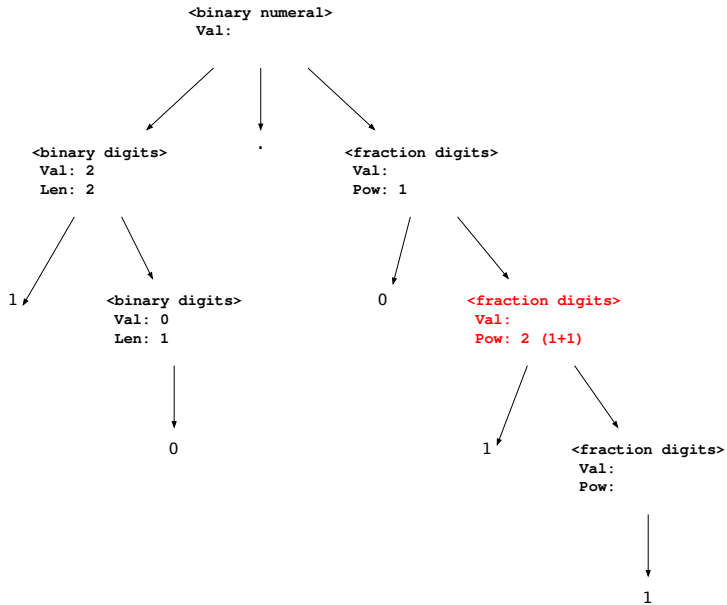


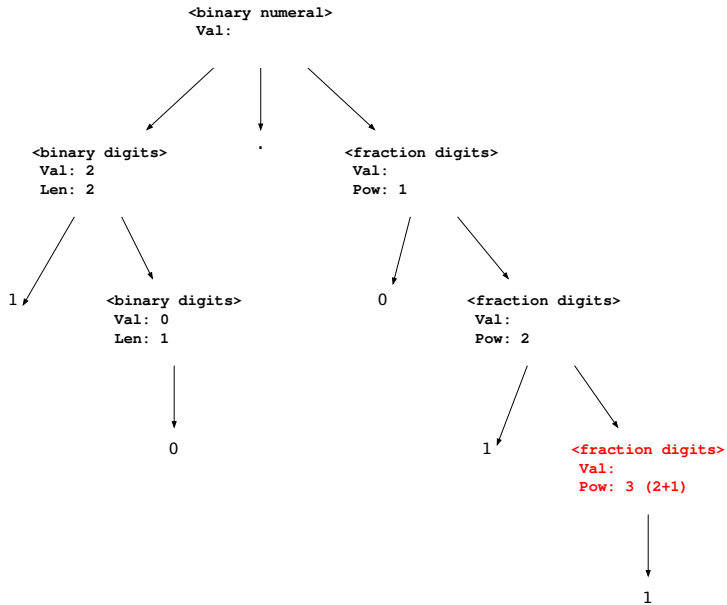


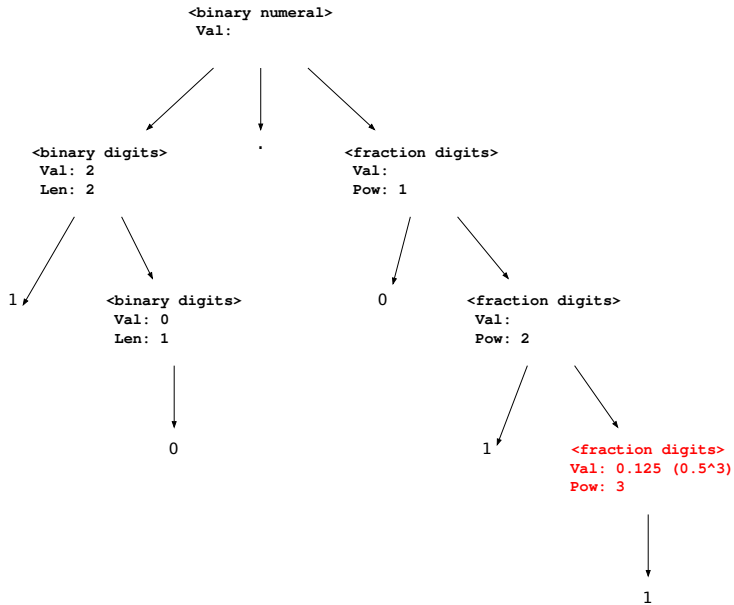


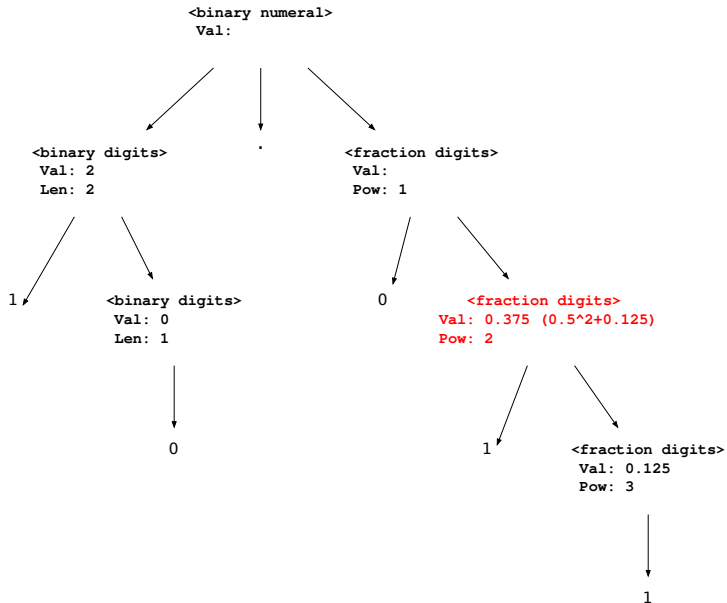


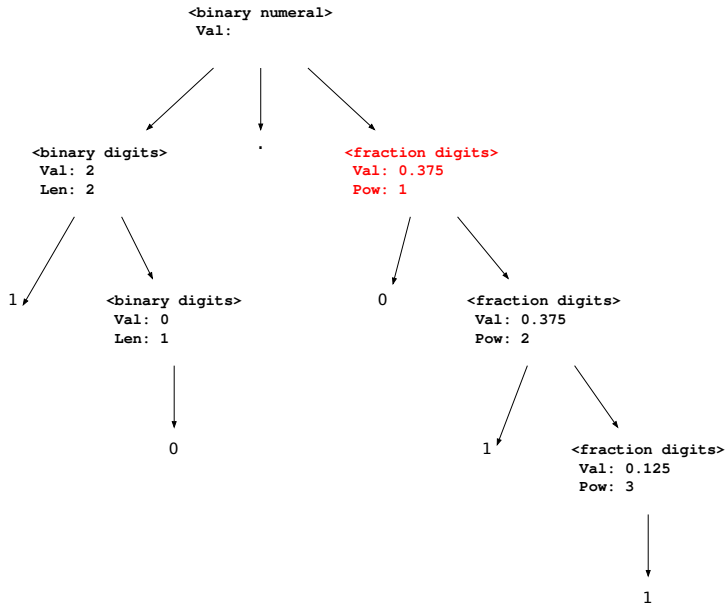




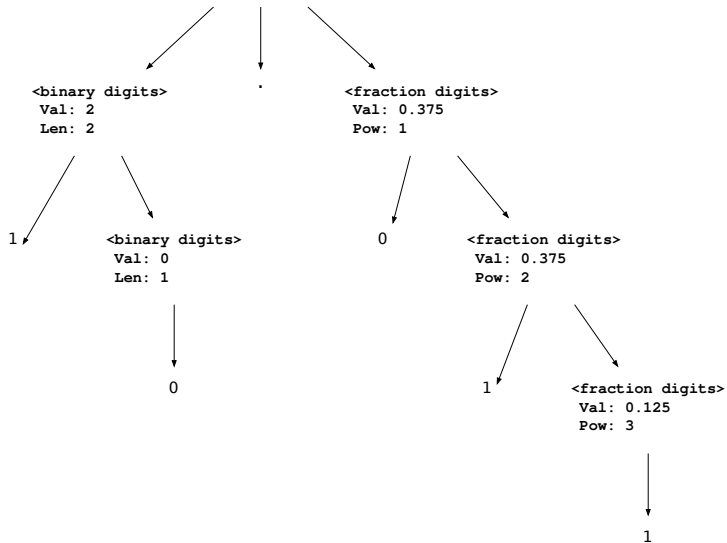




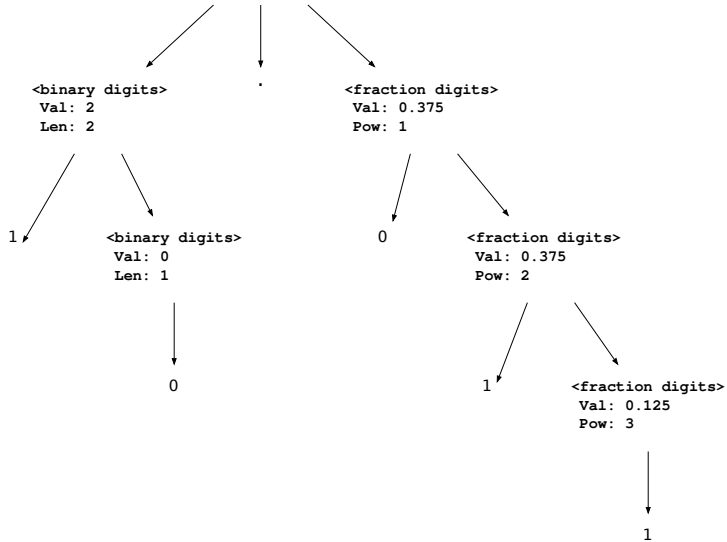




<binary numeral>
Val: 2.375 (2+0.375)



<binary numeral>
Val: 2.375



Egy jól definiált attribútum grammatikában (WAG):

- Minden lehetséges mondat minden szintaxisfájának attribútumai egyértelműen kiszámíthatóak
- Az ilyen grammatikákban a direkt attribútumfüggőségek gráfja körmentes
- Azaz biztosan létezik egy sorrend, amelyben ki tudjuk értékelni az összes attribútumot és a feltételek ellenőrizhetők
- Általánosan egy egyszerű nemdeterminisztikus algoritmussal elvégezhető (naiv megoldás)

Speciális megkötésekkel elérhető, hogy a kiszámítási sorrend könnyen meghatározható legyen (particionált, rendezett AG).

■ S-Attribútum grammatikák

- Kizárólag szintetizált attribútumokat tartalmaznak
- Parser-generátorok gyakran alkalmazzák
- Az alulról-felfelé szintaktikus elemzéshez illeszkedik

■ L-Attribútum grammatikák

- Szintetizált és örökölt attribútumokat is tartalmaz, de megkötések vannak a lehetséges függőségekre
- Az attribútumok biztosan kiszámíthatóak egyetlen bejárással
- A felülről-lefelé szintaktikus elemzésnél használják

A szintaktikus elemző generátorok olyan eszközök, amelyekkel elemzőket tudunk készíteni a nyelvhez a környezetfüggetlen grammatikája alapján.

- Tulajdonképp egy teljes lexikális és szintaktikus elemzést kapunk programozás nélkül
- Továbbá gyakran lehetővé teszik attribútumok használatát is
- Ezzel leírhatóvá válik a szemantikus elemzés és a kódgenerálás is
- Tehát fordítóprogramot készít automatikusan
- Különböző nyelveken különböző implementációk
- yacc, yecc, bisonc++, antlr, happy, ...

Az attribútum grammatikával lényegében fordítási szemantikát definiálunk.

- Mivel legtöbbször alulról-felfelé elemzőt implementálnak, S-AG-vel adják meg a szemantikát
- A kódgeneráláshoz felvesznek egy “kód” attribútumot, amely a célnyelvi szimbólumokat tartalmazza
- A kód alulról felfelé folyik a fában és végül összeépül a lefordított program
- A kezdőszimbólum kód attribútuma tartalmazza a lefordított programot

- Nyelvkiterjesztés
- Nyelvspecifikus keretrendszerek
- DSL implementációk
- Protokollok leírása
- Típusok leírása, adatgenerálás

- Környezetfüggetlen grammatikák
- Környezetfüggetlen kontra környezetfüggő
- Környezetfüggő tulajdonságok
- Attribútum grammatikák, alosztályaik, kiértékelés
- Alkalmazási területek

C++

Az $a^n b^n c^n$ környezetfüggő nyelv végrehajtható szemantikája elérhető a kurzus anyagai között. A statikus és dinamikus szemantikát is bemutató S-attribútum grammatika implementációt *flex* és *bison* C++ segítségével definiáltuk.