

## Dagens tema

- Hva er kompilering?
- Hvordan foreta syntaksanalyse av et program?
- Hvordan programmere dette i Java?
- Hvordan oppdage syntaksfeil?

# Hva er kompilering?

Anta at vi lager dette lille programmet `mini.pas` (kalt *kildekoden*):

```
/* Et minimalt Pascal-program */  
program Mini;  
begin  
    write('x');  
end.
```

Dette programmet kan ikke kjøres direkte på noen datamaskin, men det finnes en x86-kode (kalt **maskinkoden**) som gjør det samme:

```
$ od -t x1 mini
0000000 7f 45 4c 46 01 01 01 00 00 00 00 00 00 00 00
0000020 02 00 03 00 01 00 00 00 40 83 04 08 34 00 00
0000040 24 08 00 00 00 00 00 00 34 00 20 00 08 00 28
0000060 1e 00 1b 00 06 00 00 00 34 00 00 00 34 80 04
0000100 34 80 04 08 00 01 00 00 00 01 00 00 05 00 00
0000120 04 00 00 00 03 00 00 00 34 01 00 00 34 81 04
0000140 34 81 04 08 13 00 00 00 13 00 00 00 04 00 00
0000160 01 00 00 00 01 00 00 00 00 00 00 00 80 04 08
0000200 00 80 04 08 cc 05 00 00 cc 05 00 00 05 00 00
0000220 00 10 00 00 01 00 00 00 cc 05 00 00 cc 95 04
0000240 cc 95 04 08 00 01 00 00 08 01 00 00 06 00 00
```

⋮

Det er ikke lett å lese slik kode – det går bedre i **assemblerkode** som kan oversettes til maskinkode av en **assembler**:

```
# Code file created by Pascal2100 compiler 2015-09-15 22:06:52
.extern write_char
.extern write_int
.extern write_string
.globl _main
.globl main

_main:
main:    call    prog$mini_1          # Start program
        movl    $0,%eax              # Set status 0 and
        ret                                # terminate the program

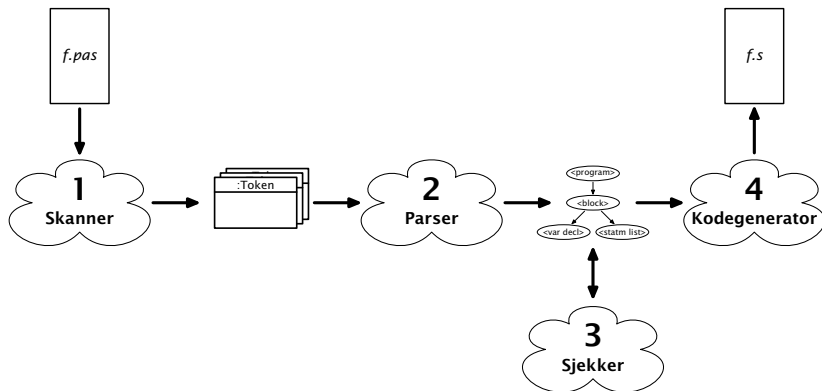
prog$mini_1:
        enter   $32,$1               # Start of mini
        movl    $120,%eax            # char 120
        pushl   %eax                 # Push param #1.
        call    write_char
        addl    $4,%esp              # Pop parameter.
        leave
        ret                          # End of mini
```

## Kompilatoren

En kompilator leser Pascal2100-koden og lager x86-assemblerkoden.

En slik kompilator skal dere lage.

## Strukturen til kompilatoren vår



## Programtreet

De færreste programmeringsspråk kan oversettes linje for linje, men det ville vært mulig med Pascal2100.

Det enkleste er likevel å lagre programmet på intern form først.

Det naturlige da er å lage et tre ved å bruke klasser, objekter og pekere. Her er OO-programmering ypperlig egnet.

Et tre er den beste representasjonsformen

# Et Pascal2100-program

Et program består av en navngitt blokk:

**program**

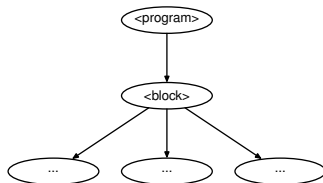




Programmet `mini.pas` representeres da av

```

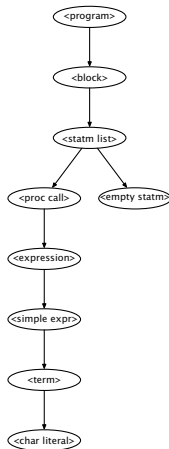
/* Et minimalt Pascal-program */
program Mini;
begin
    write('x');
end.
    
```



Det fulle programmet mini.pas ser slik ut:

```
/* Et minimalt Pascal-program */
program Mini;
begin
    write('x');
end.
```

Slike trær skal vi lage.



# Syntaksanalyse

På skolen hadde vi grammatikkanalyse hvor vi fant subjekt, predikat, indirekte og direkte objekt:

*Faren ga datteren en ball.*

(Det er ikke alltid like enkelt:

*Fanger krabber så lenge de orker. )*

Syntaksanalyse er på samme måte å finne hvilke språkelementer vi har og bygge **syntakstreet**.



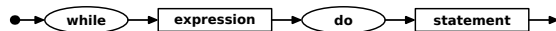
**Heldigvis:** Analyse av programmeringsspråk er enklere enn naturlige språk:

- Programmeringsspråk har en klar og entydig definisjon i jernbanediagrammer eller tilsvarende.
- Programmeringsspråk er laget for å kunne analyseres rimelig enkelt.

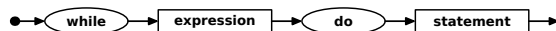
# Grammatikk

Grammatikken (i form av jernbanediagrammene) er et ypperlig utgangspunkt for å analysere et program og bygge opp syntakstreet:

## while-statm



## while-statm



Utifra dette vet vi:

- ① Først kommer symbolet *while*.
- ② Så kommer en *expression*.
- ③ Etter den kommer en *do*.
- ④ Til sist kommer et *statement*.

# Programmering av syntaksanalyse

Utifra jernbanediagrammet kan vi lage en skisse for en metode som analyserer en `while`-setning i et Pascal2100-program:

```
class WhileStatm extends Statement {  
    :  
  
    static WhileStatm parse(Scanner s) {  
        <Sjekk at vi har lest while>  
        Expression.parse(s);  
        <Sjekk at vi har lest do>  
        Statement.parse(s);  
    }  
}
```

Stort sett gjør vi to ting:

- Terminaler (i rundinger) sjekkes.
- Ikke-terminaller (i firkanter) overlates til sine egne metoder for analysering.

... og dermed har problemet nærmest løst seg selv!



## Er det så enkelt?

Mange programmeringsspråk (som Pascal men ikke Java, C og C++) er designet slik at denne teknikken kalt **recursive descent** fungerer.

Et analyseprogram for et LL(1)-språk er aldri i tvil om hvilken vei gjennom programmet som er den rette.

Ved analyse av LL(2)-språk må man av og til se ett symbol fremover. Pascal2100 er LL(2).<sup>1</sup>

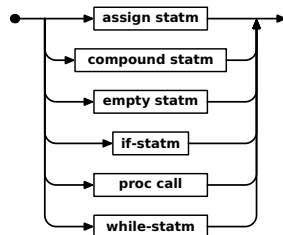


<sup>1</sup>Det finnes ett unntak; det skal vi se på neste uke.

Veien er alltid gitt!

# Veien er alltid klar

## statement



# Hvordan programmerer vi dette i Java?

Husk at målet med analysen er tofoldig:

- ① Vi skal sjekke at programmet er riktig.
- ② Vi skal bygge opp et korrekt syntakstre.

Det er naturlig å koble analysemetoden til den klassen som skal inngå i syntakstreet.

- Hver ikke-terminal i jernbanediagrammet implementeres av en Java-klasse.<sup>2</sup>
- Hver av disse klassene får en metode

```
static xxx parse(Scanner s) { ... }
```

som kan analysere «seg selv» og returnere et objekt som representerer «seg selv».

## Hvordan skrive dette i Java?

Pascal2100-kompilatoren *kan* ha disse klassene til å representere **ikke-terminaler**. (Abstrakte klasser er satt i hakeparenteser og innrykk markerer subklasser.)

CodeFile	- [Operator]	- - CompoundStatm
LogFile	- - FactorOperator	- - EmptyStatm
Main	- - PrefixOperator	- - IfStatm
[PascalSyntax]	- - RelOperator	- - ProcCallStatm
- Block	- - TermOperator	- - WhileStatm
- - Library	- ParamDeclList	- StatmList
- ConstDeclPart	- [PascalDecl]	- Term
- Expression	- - ConstDecl	- [Type]
- [Factor]	- - EnumLiteral	- - ArrayType
- - [Constant]	- - ParamDecl	- - EnumType
- - - CharLiteral	- - ProcDecl	- - RangeType
- - - NamedConst	- - - FuncDecl	- - TypeName
- - - NumberLiteral	- - Program	- TypeDeclPart
- - - StringLiteral	- - TypeDecl	- VarDeclPart
- - FuncCall	- - VarDecl	Scanner
- - InnerExpr	- SimpleExpr	Token
- - Negation	- [Statement]	TokenKind
- - Variable	- - AssignStatm	

# Statiske deklarasjoner

## Vanlige variabler i klasser

Vanlige variabler oppstår når et objekt opprettes. Det kan derfor være vilkårlig mange av dem.

## static-variabler

Disse ligger i «selve klassen» så det vil alltid være nøyaktig én av dem.

## Et eksempel

```
class Item {
    private static int total = 0;
    public int id;

    public Item() { id = ++total; }
}

class RunItem {
    public static void main(String arg[]) {
        Item a = new Item(), b = new Item();

        System.out.println("a.id = "+a.id);
        System.out.println("b.id = "+b.id);
    }
}
```

## Resultatet

```
a.id = 1
b.id = 2
```

## Vanlige metoder i klasser

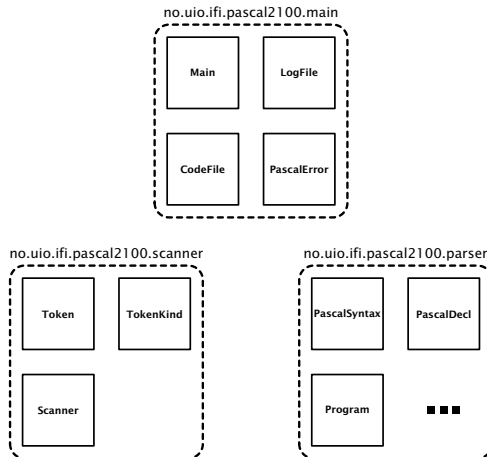
Vanlige metoder ligger logisk sett i det enkelte objektet. Når de refererer til variabler, menes variabler i samme objekt eller `static`-variabler i klassen.

## `static`-metoder

Disse ligger logisk sett i «selve klassen». De kan derfor kalles før noen objekter er opprettet men de kan bare referere til `static`-variabler.



# Prosjektet vårt i Java



# Rot-klassen

Vi lar alle parser-klassene være subklasser av PascalSyntax.

```
package no.uio.ifi.pascal2100.parser;

import no.uio.ifi.pascal2100.main.*;

public abstract class PascalSyntax {
    public int lineNumber;

    PascalSyntax(int n) {
        lineNumber = n;
    }

    boolean isInLibrary() {
        return lineNumber < 0;
    }

    abstract void check(Block curScope, Library lib);
    abstract void genCode(CodeFile f);
    abstract public String identify();
    abstract void prettyPrint();

    void error(String message) {
        Main.error("Error at line " + lineNumber + ": " + message);
    }
}
```



Alle deklarasjonene er også subklasser av `PascalDecl`:

```
package no.uio.ifi.pascal2100.parser;

public abstract class PascalDecl extends PascalSyntax {
    String name, progProcFuncName;
    int declLevel = 0, declOffset = 0;
    Type type = null;

    PascalDecl(String id, int lNum) {
        super(lNum);
        name = id;
    }
}
```

## Klassen Program kan da se slik ut:

```
package no.uio.ifi.pascal2100.parser;

import no.uio.ifi.pascal2100.main.*;
import no.uio.ifi.pascal2100.scanner.*;
import static no.uio.ifi.pascal2100.scanner.TokenKind.*;

/* <program> ::= 'program' <name> ';' <block> '.' */

public class Program extends PascalDecl {
    Block progBlock;

    Program(String id, int lNum) {
        super(id, lNum);
    }

    @Override public String identify() {
        return "<program> " + name + " on line " + lineNum;
    }
}
```

```

public static Program parse(Scanner s) {
    enterParser("program");
    s.skip(programToken);
    s.test(nameToken);

    Program p = new Program(s.curToken.id, s.curLineNum());
    s.readNextToken();
    s.skip(semicolonToken);
    p.progBlock = Block.parse(s); p.progBlock.context = p;
    s.skip(dotToken);

    leaveParser("program");
    return p;
}
    
```

# Husk

I scanner.Scanner har vi

```
public void test(TokenKind t) {  
    if (curToken.kind != t)  
        testError(t.toString());  
}  
  
public void testError(String message) {  
    Main.error(curLineNum(),  
               "Expected a " + message +  
               " but found a " + curToken.kind + "!");  
}  
  
public void skip(TokenKind t) {  
    test(t);  
    readNextToken();  
}
```

og i main.Main har vi

```
public static void error(String message) {  
    log.noteError(message);  
    throw new PascalError(message);  
}
```



## Samarbeid med skanneren

Hvordan sikrer vi at symbolstrømmen fra Scanner er i fase med vår parsering?

Det beste er å vedta noen regler som alle parse-metodene *må* følge:

- ➊ Når man kaller parse, skal første symbol være lest inn!
- ➋ Når man returnerer fra en parse, skal første symbol *etter* konstruksjonen være lest!

Hvordan finner man programmeringsfeil?

## Feilsjekking

Sjekken på syntaksfeil er svært enkel:

### Hvordan finne feil?

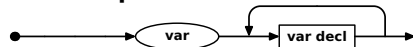
Hvis neste symbol ikke gir noen lovlig vei i diagrammet, er det en feil.



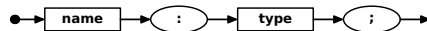
## Hvordan finner man programmeringsfeil?

### Et eksempel

#### var decl part



#### var decl



var v : integer = 5 ;

## Parsing av While-setning

```
package no.uio.ifi.pascal2100.parser;

import no.uio.ifi.pascal2100.main.*;
import no.uio.ifi.pascal2100.scanner.*;
import static no.uio.ifi.pascal2100.scanner.TokenKind.*;

/* <while-statm> ::= 'while' <expression> 'do' <statement> */

class WhileStatm extends Statement {
    Expression expr;
    Statement body;

    WhileStatm(int lNum) {
        super(lNum);
    }

    @Override public String identify() {
        return "<while-statm> on line " + lineNum;
    }

    static WhileStatm parse(Scanner s) {
        enterParser("while-statm");

        WhileStatm ws = new WhileStatm(s.curLineNum());
        s.skip(whileToken);

        ws.expr = Expression.parse(s);
        s.skip(doToken);
        ws.body = Statement.parse(s);

        leaveParser("while-statm");
        return ws;
    }
}
```

