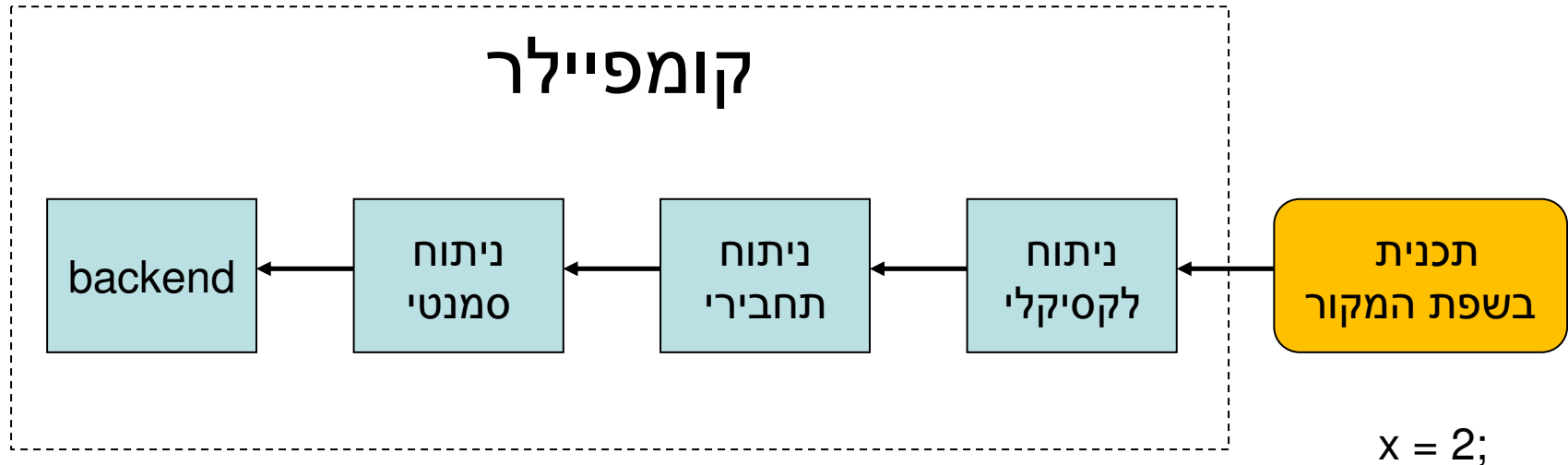


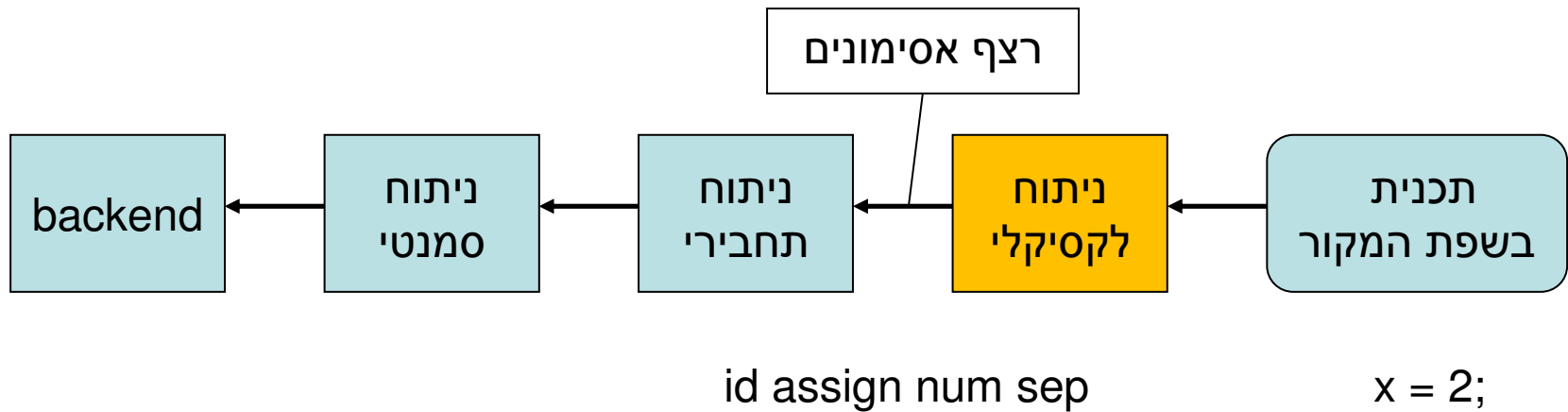
ניתוח לקסיקלי וכלי Lex

נכתב ע"י אלכס קוגן (sakogan@cs)
סמסטר חורף, תשס"ח

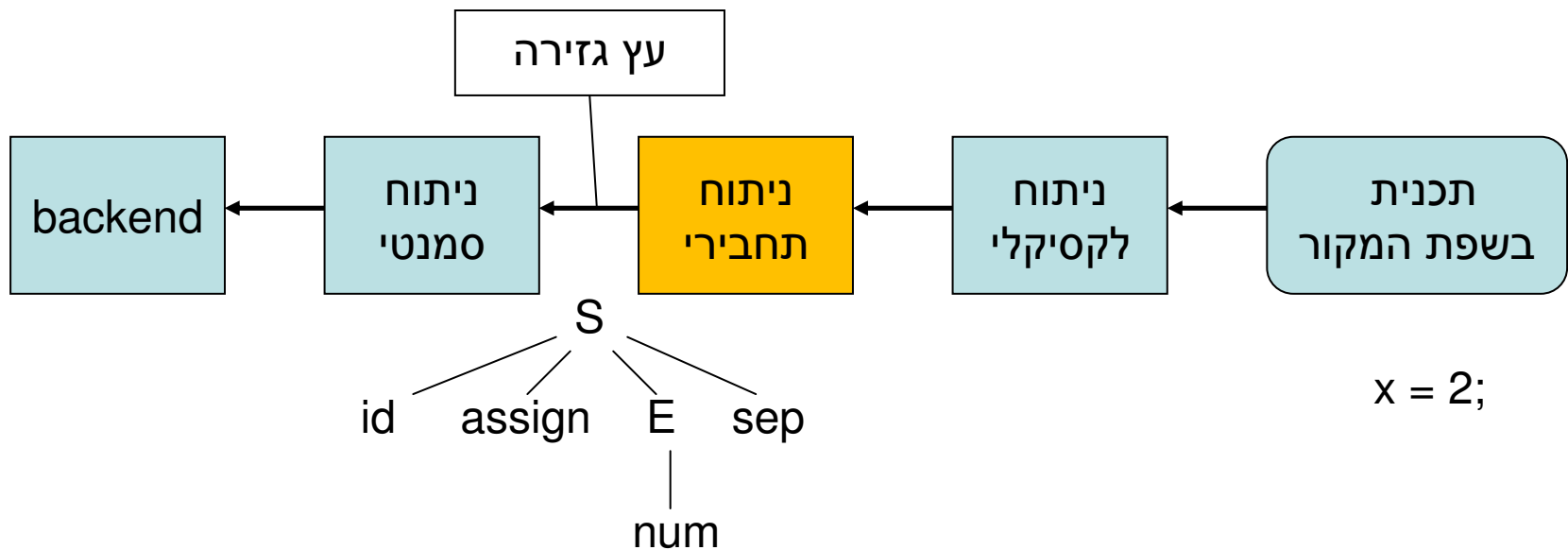
מבנה הקומפילר



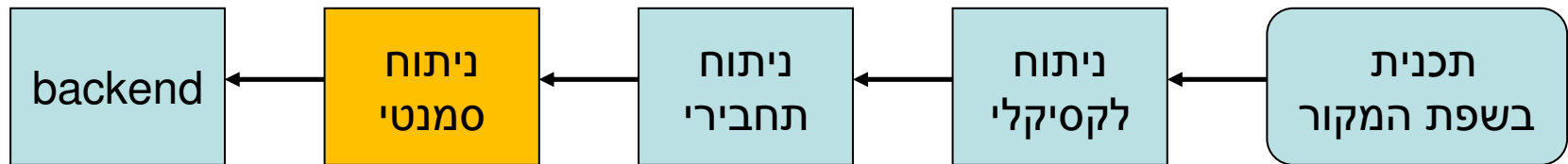
מבנה הקומפילר



מבנה הקומפילר



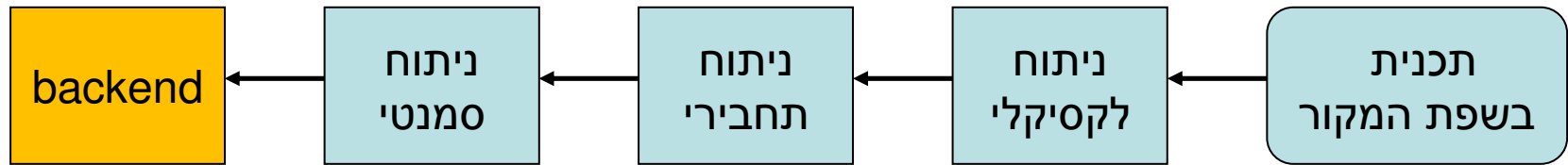
מבנה הקומפילר



האם קיים משתנה ששמו x?
האם x מטיפוס int?

x = 2;

מבנה הקומפילר



ייצור קוד בשפת היעד

$x = 2;$

הגדרות

- **אסימון (token):** יחידה בסיסית המשמשת כטרמינל בדקדוק שגוזר את שפת התכנות
- **לקסמה (lexeme):** מחרוזת בקלט שהמנתח הלקסיקלי התאים לאסימון כלשהו
- **תבנית (pattern, regexp):** ביטוי רגולרי שמגדיר את ההתאמה בין אוסף הלקסמות לאסימון מסוים – דוגמה: ניתן להגדיר את האסימון `num` ע"י התבנית $(0+1+\dots+9)^+$

תפקיד המנתח הלקסיקלי

- מחלק את קוד המקור לאסימונים ("מילים")
- מסנן חלקים שאינם דרושים להמשך הניתוח
 - למשל, רווחים, ירידות שורה
 - מה עם comments?
- מזהה שגיאות לקסיקליות - מחרוזות שאינן להיות אף אסימון
 - למשל, "@" בשפת C

המנתח הלקסיקלי - דוגמה

- ניתוח לקסיקלי עבור "x = y + 2":

לקסמות	x	=	y	+	2	\n
אסימונים	id	assign	id	op	num	

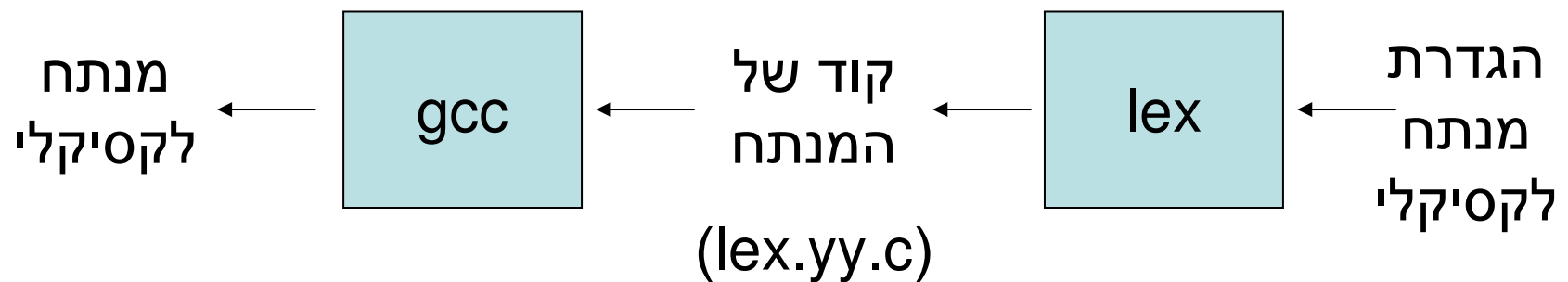
המנתח לא
יחזיר אסימון

המנתח הלקסיקלי ותכונות סמנטיות

- הלקסמות אינן עוברות הלאה לשלבים הבאים בניתוח
 - הנמתח התחבירי אינו צריך אותן
 - המנתח הסמנטי כן צריך מידע עליהן:
 - האם קיים משתנה ששמו x?
 - איזו פעולה בדיוק הופעלה?
- המנתח הלקסיקלי מחשב לכל אסימון תכונות סמנטיות:
 - `id {name = "x"}`
 - `op {type = "+"}`
- התכונות הסמנטיות ידונו בהרחבה בהמשך

כלי Lex

- כלי לייצור מנתחים לקסיקליים
- צורת הפעלה:



- בקורס נעבוד עם כלי תואם Flex

בניית מנתח לקסיקלי בעזרת Lex

- עריכת קובץ הגדרת המנתח (source.lex)
- flex source.lex
- gcc -ll lex.yy.c

מבנה קובץ הגדרות המנתח

- קובץ הגדרות מורכב משלושה חלקים המופרדים בעזרת שורות שמכילות "%%" בלבד:

Definition section

%%

Rules section

%%

C code section

דוגמה למבנה קובץ הגדרות המנתח

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace     ([\t\n ])
12
13 %%
14 {digit}+       showToken("number");
15 {letter}+@{letter}+\.com showToken("email address");
16 {whitespace}  ;
17 .             printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```

דוגמה למבנה קובץ הגדרות המנתח

Definitions
section

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace     ([\t\n ])
12
13 %%
14 (digit)+       showToken("number");
15 (letter)+@{(letter)+\.com showToken("email address");
16 (whitespace)   ;
17 .              printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```

דוגמה למבנה קובץ הגדרות המנתח

הגדרות של שפת C -
קוד מועתק כפי שכתוב
לקובץ ש-Lex מייצר

אופציות השולטות על
צורת העבודה של Flex

הגדרת מקרואים
בעזרת ביטויים
רגולריים - לשימוש
בחלק הבא

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace     ([\t\n ])
12
13 %%
14 (digit)+       showToken("number");
15 (letter)+@{(letter)+\.com showToken("email address");
16 (whitespace)   ;
17 .              printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```


ביטויים רגולריים של Lex

משמעות	ביטוי רגולרי "רגיל"	ביטוי רגולרי של Lex
התו a	a	a
כל תו פרט לירידת שורה	$\Sigma \setminus \{\backslash n\}$.
אחד מהתווים שבתוך הסוגריים	x+y+z a+b+c+...+z	[xyz] [a-z]
מספר כלשהו של z-ים כולל	r*	r*
אפס / לא כולל אפס	r+	r+

דוגמאות נוספות ניתן למצוא באתר הקורס

דוגמה למבנה קובץ הגדרות המנתח

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace    ([\t\n ])
12
13 %%
14 {digit}+      showToken("number");
15 {letter}+@{letter}+\.com showToken("email address");
16 {whitespace} ;
17 .            printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```

Rules section
הגדרות של
אסימונים ופעולות
ש-Lex צריך לבצע
בעת זיהוי שלהם

דוגמה למבנה קובץ הגדרות המנתח

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace     ([\t\n ])
12
13 %%
14 {digit}+       showToken("number");
15 {letter}+@{letter}+\.com showToken("email address");
16 {whitespace}  ;
17 .             printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```

C code section
הגדרות פונקציות
שהוכרוזו בחלק ראשון

משתנים גלובליים
של Lex

משתנים גלובליים של Lex

שם	טיפוס	משמעות
<code>yytext</code>	<code>char *</code>	הטקסט של הלקסמה האחרונה שזוהתה
<code>yy leng</code>	<code>int</code>	אורך הלקסמה האחרונה שזוהתה
<code>yylineno</code>	<code>int</code>	השורה הנוכחית בקלט
<code>yyval</code>	<code>user defined</code>	משמש לתקשורת עם המנתח התחבירי (פרטים בהמשך...)

דוגמה למבנה קובץ הגדרות המנתח

```
1  %{
2  /* Declarations section */
3  #include <stdio.h>
4  void showToken(char *);
5  %}
6
7  %option yylineno
8  %option noyywrap
9  digit          ([0-9])
10 letter         ([a-zA-Z])
11 whitespace     ([\t\n ])
12
13 %%
14 {digit}+       showToken("number");
15 {letter}+@{letter}+\.com showToken("email address");
16 {whitespace}  ;
17 .             printf("Lex doesn't know what that is!\n");
18 %%
19
20 void showToken(char * name) {
21     printf("Lex found token %s , name);
22     printf("the lexeme is %s , yytext);
23     printf("its length is %d\n", yyleng);
24 }
```

תכונות המנתח הלקסיקלי ש-Lex בונה

- קורא קלט מ-stdin וכותב פלט ל-stdout
- רוב העבודה נעשית בפונקציה yylex שתפקידה:
 - לקרוא את הקלט
 - לזהות אסימונים
 - לבצע את פעולה המתאימה לאסימון
- yylex חוזרת רק כאשר:
 - המשתמש כתב return בפעולה של אסימון
 - מגיעה לסוף קלט

טיפול בשגיאות ניתוח

- שגיאת ניתוח נוצרת כאשר בנקודה מסוימת yylex לא מצליחה לזהות אף אסימון
- במקרה כזה, היא מעתיקה את התו הנוכחי לפלט ומנסה שוב מהתו הבא
- דוגמה:

abc\$178

מזוהה כ-id מזוהה כ-num

לא מזוהה ולכן מודפס לפלט

טיפול בשגיאות ניתוח (המשך)

- כדי לתת הודעת שגיאה יותר מפורטת, ניתן להשתמש באסימון הנקודה (.) בתור אסימון אחרון

– כפי שנעשה בדוגמה:

```
13  %%  
14  {digit}+           showToken("number");  
15  {letter}+@{letter}+\.com showToken("email address");  
16  {whitespace}      ;  
17  .                  printf("Lex doesn't know what that is!\n");  
18  %%
```


פתרון קונפליקטים

- קונפליקט נוצר כאשר אותו קלט יכול להתאים למספר אסימונים
- דוגמה: נניח שמוגדרים האסימונים הבאים:

for "for"

id [a-z]+

מחרוזות שיגרמו לקונפליקט	ניתוחים אפשריים
abc	
ford	
for	

פתרון קונפליקטים

- קונפליקט נוצר כאשר אותו קלט יכול להתאים למספר אסימונים
- דוגמה: נניח שמוגדרים האסימונים הבאים:

for “for”

id [a-z]+

מחרוזות שיגרמו לקונפליקט	ניתוחים אפשריים
abc	a, ab, abc
ford	for, ford, f, fo
for	for (id) / for (for)

הכללים לפתרון הקונפליקט

1. המנתח חמדן – תמיד מעדיף את הלקסמה
הארוכה ביותר שניתן לבחור
2. אם כלל (1) לא פתר את הקונפליקט, בוחרים
באסימון בעל עדיפות גבוהה יותר
 - האסימון שמופיע ראשון

פתרון הדוגמה

- נניח שמוגדרים האסימונים הבאים:

for “for”

id [a-z]+

לפי כלל	הניתוח שייבחר	ניתוחים אפשריים	מחרוזות שיגרמו לקונפליקט
1	abc	a, ab, abc	abc
1	ford	for, ford, f, fo	ford
2	for (for)	for (id) / for (for)	for

מקורות נוספים

- דוגמאות לשימוש ב-Lex ניתן למצוא באתר הקורס
- שימוש ב-manual מותקן על Lex וביטויים רגולריים:
man lex –
man –s 5 regexp –
- אתר הבית של כלי Flex:
<http://flex.sourceforge.net>
- מדריך לשימוש ב-Flex:
http://www.gnu.org/software/flex/manual/html_node/flex_toc.html