

עבודת גמר

לקבלת תואר טכנאי תוכנה

נושא: The Easy Compiler

מהדר לשפת Easy



The Easy Compiler



המגיש: אביב אש

ת.ז המגיש: 214887556

שם המנחה: מיכאל צ'רנובילסקי

תשפ"ה

אפריל 2025

תוכן עניינים

2.....	תוכן עניינים
4.....	תקציר
5.....	מושגים
5.....	מושגים כללים
6.....	תיאור הפרויקט
6.....	תכולת השפה
6.....	תיאור אבני השפה
6.....	אופרטורים בינאריים:
7.....	משתנים בסיסיים:
8.....	השמת נתונים
8.....	השמה מורכבת
9.....	השוואות
9.....	תנאים לוגיים
10.....	תנאים
11.....	פונקציות
11.....	קריאה לפונקציה
12.....	מצביעים
13.....	תוכנית לדוגמא בשפת easy
14.....	דקדוק השפה
17.....	רקע תאורטי בתחום הפרויקט
17.....	קומפיילר/מהדר – הגדרה כללית:
19.....	סוגים שונים של מהדרים:
19.....	סוגים שונים של שפות תכנות:
19.....	שלבי הקומפילציה:
19.....	Front end:
19.....	ניתוח לקסיקלי (Lexical Analysis)
20.....	ניתוח תחבירי (Syntax Analysis)
20.....	דקדוק חסר הקשר (Context free grammar)
22.....	סוגי מנתחים תחביריים
22.....	מלמעלה למטה (top - down)
23.....	מנתחי bottom-up
23.....	מלמטה למעלה (bottom - up)
24.....	ניתוח סמנטי (Semantic Analysis) –
24.....	-Back end
24.....	אופטימיזציה –

24.....	יצירת קוד הסף –
24.....	טיפול בשגיאות-.....
25.....	תיאור הבעיה האלגוריתמית.....
26.....	אסטרטגיה נבחרת לפתרון.....
26.....	האסטרטגיה לפיתוח המנתח הסמנטי.....
27.....	האוטומט של tec
30.....	האסטרטגיה לפיתוח המנתח התחבירי.....
30.....	הפעולות האפשריות
33.....	האסטרטגיה לפיתוח המנתח הסמנטי.....
33.....	האסטרטגיה לפיתוח שלב כתיבת קוד הסף.....
36.....	ארכיטקטורה של הפתרון המוצע בפורמט של Top-Down Level Design.....
37.....	תרשים Top-down מפורט.....
38.....	מבני נתונים.....
38.....	מטריצה.....
39.....	מפת hash.....
39.....	מחסנית.....
40.....	עץ.....
40.....	מערך.....
41.....	תיאור סביבת העבודה ושפות התכנות.....
42.....	אלגוריתם ראשי.....
42.....	אלגוריתם העל
42.....	אלגוריתם המנתח הלקסיקלי
42.....	מציאת אסימון הבא.....
43.....	אלגוריתם המנתח התחבירי
43.....	Shift.....
43.....	Reduce.....
44.....	אלגוריתם המנתח הסמנטי.....

תקציר

בחרתי לעשות את פרויקט הגמר שלי בנושא פיתוח מהדר (compiler)

כבר בתחילת לימודי במדעי המחשב מצאתי את המהדר כרכיב מסתורי וקסום, איך זה ייתכן שאני כותב קוד באנגלית והמחשב יודע לפרש זאת ולהריץ אותו, זאת אחת הסיבות שבגללן בחרתי את נושא זה כי איזו דרך יותר טובה להבין משהו מללכלך את הידיים ולעשות אותו בעצמך.

במהלך הפרויקט חקרתי מספר כלי הקומפילציה ותרגום (אינטרפטציה) שונים ביניהם (cpython,gcc) בשביל להבין יותר על הנושא וקריאה יותר עם הנושאים האלו עניינה אותי מאוד ועזרה לי בכתיבת והבנת הפרויקט

הפרויקט פיתח אותי מאוד אני מרגיש שהשתפרתי מאוד כמתכנת. למדתי לכתוב קוד יותר קריא ודינמי וזהו יתרון גדול. במהלך הכתיבה הבנתי שהרבה פעמים ארצה לשנות או להוסיף פונקציונליות ברכיב בקוד שלי ואם הקוד לא כתוב כראוי. למדתי שפת תכנות חדשה, החלטתי לכתוב את הפרויקט שלי בשפת ++c בשל גמישות השפה והיותה גם מונחת עצמים וגם זמן הריצה המהיר שלה וגישה לזיכרון, תהליך למידת השפה היה משמעותי ומקדם עבורי.

בקיצור זהו הפרויקט הגדול והמורכב ביותר שלי עד כה, פגשתי אתגרים רבים ולמדתי המון.

מקווה שתהנו 😊

מושגים

בפרק זה אציג את המושגים השונים הקשורים ביצירת מהדר

מושגים כלליים

מהדר (compiler) - רכיב תוכנה שמטרתו לתרגם קוד מקור (Source Code) לשפת יעד, כמו קובץ הרצה (Executable) או שפת ביניים (Intermediate Representation - IR). מהדרים נפוצים: gcc, clang

אסימון (token) – אסימון הוא חלק בעל משמעות בקוד, לכל אסימון יש סוג ושם. לדוגמא בשורת הקוד while(x) נמצאים שלושה אסימונים:

[while : **keyWord**] , [(**openParen**)] , [**x**: **identifier**] , [**closedParen**])

עץ ניתוח תחבירי (AST) – עץ המורכב מאסימונים המתאר את המבנה הסמנטי של התוכנית

טבלת הסימנים (Symbol table) - מבנה נתונים שמנהל מידע על משתנים, פונקציות, וטיפוסים בקוד.

דקדוק חופשי-הקשר (CFG) הינו חלק מענף במתמטיקה הנקרא הוא חלק מענף המתמטיקה שנקרא **תורת השפות הפורמליות** (Formal language Theory) בעזרת דקדוק זה נבטא את המבנה של התוכנית (הסבר מפורט יותר בחלק **הרקע התאורטי**)

תיאור הפרויקט

תכולת השפה

השפה easy שפה יחסית מינימליסטית והיא נבנתה כשמה, היא אמורה להיות קלה ללמידה ואינטואיטיבית לאנשים עם רקע בתכנות ולוקחת הרבה השראה מ C למרות שהיא כוללת גם מאפיינים משפות אחרות כמו Java ו JS

היא שפה פרוצדורלית ולכן מכילה את המבני השפה הנחוצים לה

השפה מבנים בסיסיים לולאות, תנאים, והגדרת ביטויים מתמטיים אך היא מכילה מבנים נוספים כגון: פונקציות, מצביעים ומערכים

תיאור אבני השפה

אופרטורים בינאריים:

כמו בהרבה שפות גם easy מכילה מגוון אופרטורים בינאריים

אופרטור בינארי הינו אופרטור המפעיל פעולה בין שני ביטויים

<ביטוי 1> <אופרטור> <ביטוי 2>

האופרטורים:

+ חיבור

- חיסור

* כפל

/ חילוק

~ פעולת ביטים not

| פעולת ביטים or

& פעולת ביטים and

^ פעולת ביטים xor

משתנים בסיסיים:

השפה כוללת מספר סוגי משתנים בסיסיים: `int`, `char`, `float`:

הגדרת משתנה תתבצע בצורה הבאה (בדומה לC):

<ערך אתחול ראשוני> = <שם המשתנה> <סוג המשתנה>

או:

<שם המשתנה> <סוג המשתנה>;

יש לשים לב כי בדרך האתחול השנייה המשתנה יאותחל לערך "זבל" ויכלול את תכולת הזכרון הקודם במקום שהוא תפס מבלי לדרוס אותו.

דוגמאות לאתחול נתונים בשפה

```
int x1 = 5;

char y1 = 6;

int x2 = x1*2 + 2;

float y2 = 6.0;

int c;
```

השמת נתונים

ההשמה מתבצעת בדומה לאתחול רק ללא שם המשתנה בתחילת ההשמה
צד ימין של ההשמה לא יכול להיות ריק

<ערך אתחול ראשוני> = <שם המשתנה>

```
x1 = 5;
```

דוגמה להשמה תקינה:

```
x1;
```

דוגמה להשמה לא תקינה:

השמה מורכבת

השמה מורכבת הינה השמה שבנוסף להשמה מתבצעת גם פעולה נוספת

<ערך אתחול ראשוני> <אופרטור השמה מורכב> <שם המשתנה>

לדוגמא:

```
y1 += 6;
```

הביטוי לא ישים את הערך 6 במשתנה y1 אלא יוסיף אילו 6 המהדר יפרש ביטוי זה כ –

```
y1 = y1 + 6;
```

להלן כל האופרטורים להשמה מורכבת שהשפה כוללת:

- += מוסיף למשתנה בצד ימין את הביטוי בצד שמאל
- -= מחסר למשתנה בצד ימין את הביטוי בצד שמאל
- /= מחלק למשתנה בצד ימין את הביטוי בצד שמאל
- *= כופל למשתנה בצד ימין את הביטוי בצד שמאל
- ~= מבצע פעולת סיביות NOT למשתנה בצד ימין את הביטוי בצד שמאל
- |= מבצע פעולת סיביות OR למשתנה בצד ימין את הביטוי בצד שמאל
- &= מבצע פעולת סיביות AND למשתנה בצד ימין את הביטוי בצד שמאל
- ^= מבצע פעולת סיביות XOR למשתנה בצד ימין את הביטוי בצד שמאל

השוואות

השוואה מחזירה אחד משתי ערכים **0** (שקר) או **1** (אמת)
השוואה מתבצעת בעזרת האופרטורים **==** (שווה) או **!=** (לא שווה)

דוגמא 1:

$$x \% 3 == 0$$

תחזיר **אמת** אם המספר מתחלק בשלוש ללא שארית אחרת **שקר**

דוגמא 2:

$$3 != 0$$

תחזיר **שקר** כי **0** הוא לא **3!**

תנאים לוגיים

בין שני ביטויים לניתן לשים אופרטור לוגי בצורה הבאה
<ביטוי 1> && <ביטוי 2>

או

<ביטוי 1> || <ביטוי 2>

&& - מסמן "ו-" לוגי והוא יחזיר 1 (אמת) רק אם שני הביטויים שונים מאפס (אמת)
&& - מסמן "או" לוגי והוא יחזיר 1 (אמת) אם אחד או יותר משני הביטויים שונים מאפס (אמת)

לדוגמא:

$$x \% 3 == 0 \&\& x \% 2 == 0$$

יחזיר האם המספר מתחלק בשלוש ללא שארית וגם מתחלק בשניים ללא שארית

תנאים

תנאים בשפה יכתבו בצורה הבאה:

if(<Expression>) <body>

אפשר להוסיף

Else <body>

כאשר תוצאת הביטוי מניבה מספר

<Body> הינו או ביטוי אחר או מספר ביטויים בתוך סוגריים מסולסלות

אם המספר לא שווה אפס התוכנית תמשיך למה שקיים בתוך גוף התנאי

כשערכו של המספר הוא אפס התוכנית לא תכנס לגוף התנאי ותכנס לגוף **Else** (אם יש)

לדוגמא:

```
if(x) {x = 5;}
```

עם x שונה מ-0 ערכו ישתנה ל-5 אחרת - ישאר אותו הדבר

פונקציות

מהי פונקציה? פונקציה היא מבנה בשפת תכנות המאפשר לבצע קטע קוד מסוים באופן חוזר תוך מתן אפשרות להעברת פרמטרים והחזרת ערך.

פונקציות בשפה easy נכתבות כך –

{<גוף הפונקציה> <סוג משתנה> => (<רשימת פרמטרים>) <שם הפונקציה>

גוף הפונקציה הינו אפס או יותר ביטויים אשר אחד מהביטויים הינו **ret** ביטוי

השפה מכילה את מילת המפתח **ret** כאשר התוכנית מגיעה אל המילה הזאת היא תצא ותחזיר את הערך אשר נמצא אחריה

דוגמאות לפונקציות:

דוגמא 1:

פונקציה שלא מקבלת כלום ומחזירה את הערך 2

```
func1() => int {
    ret 2;
}
```

דוגמא 2:

פונקציה המקבלת מספר ומחזירה את החצי שלו

```
func2(int x) => int {
    ret x/2;
}
```

קריאה לפונקציה

קריאה לפונקציה תראה כך

(<פרמטרים>) <שם הפונקציה>

לדוגמא:

```
func2(6)
```

הפונקציה func2 תקרא עם הערך 6 ותחזיר את הערך 2

מצביעים

עוד בונוס קטן לשפה, **מצביעים ומערכים**.

בשביל שמירה על הפשטות בשפה easy המצביעים מגיעים עד לרמה אחת של ריחוק
השפה לא מאפשרת **מצביע למצביע** או כל מבצע ברמה גבוהה יותר מהצבעה לערך
מצביעים כתובים כך:

הגדרת מצביע:

< ערך אתחול ראשוני > = <סוג משתנה> * <שם משתנה>

הגדרת מערך:

< רשימת אתחול > = <סוג משתנה> [<מספר איברים במערך>] <שם משתנה>

גזירת מצביעה (dereferencing)

<שם מצביע>&

דוגמא לשימוש:

```
int[2] arr = {1,2};  
int *ptr = &arr;
```

תוכנית לדוגמא בשפת easy

להלן תוכנית המחזירה את המספר הגדול ביותר :

```
fn main() => int
{
    int x = 7*7*7;
    int y = 6;
    print(max(x,y));

    ret 0;
}

fn max(int x,int y) => int
{
    int res;

    if(x > y) res = x;
    if(y > x) res = y;

    ret res;
}
```

דקדוק השפה

מה הוא דקדוק?

דקדוק הינו אוסף של חוקים המגדירים את השפה.

להלן הדקדוק מוצג בBNF אשר הינו פורמט סטנדרטי לייצוג שפות

Program ::= <FunctionDecl> | ε

"FunctionDecl" ::= "fn" IDENTIFIER "(" <ParamList> ")" "-" <Type> "{" <StmtList> ">

ParamList ::= <Param> <ParamTail> | ε

ParamTail ::= "," <Param> <ParamTail> | ε

Param ::= <Type> IDENTIFIER

<Type> ::= <BaseType> <TypeTail>

"BaseType" ::= "int" | "float" | "char"

TypeTail ::= "[" INTEGER_LITERAL "]" | "*" | ε

StmtList ::= <Stmt> <StmtList> | ε

Stmt ::= <VarDeclStmt> | <AssignStmt> | <IfStmt> | <WhileStmt> | <ForStmt> | <>
<<ReturnStmt> | <ExprStmt>

<Body> ::= "{" <StmtList> "}" | <Stmt>

"," <VarDeclStmt> ::= <VarDeclExpr>

<VarDeclExpr> ::= <Type> IDENTIFIER <InitOpt>

InitOpt ::= "=" <AssignValue> | ε

"AssignValue> ::= <Expr> | "{" <ExprList> ">

"," <AssignStmt> ::= <AssignExpr>

<AssignExpr> ::= <AssignTarget> <AssignOp> <Expr>

AssignTarget> ::= IDENTIFIER | IDENTIFIER "[" <Expr> "]" | "*" IDENTIFIER>

"=~" | "=^" | "=|" | "=&" | "=*" | "=/" | "=-" | "=+" | "=" ::= <AssignOp>

IfStmt> ::= "if" "(" <ConditionOp> ")" <Body> | "if" "(" <ConditionOp> ")" <Body> "else>

<WhileStmt> ::= "while" "(" <ConditionOp> ")" <Body>

<ConditionOp> ::= <Expr> | <AssignExpr>

<ForStmt> ::= "for" "(" <ForInit> ";" <ExprOpt> ";" <ForUpdate> ")" <Body>

ForInit> ::= <VarDeclExpr> | <AssignExpr> | ε>

ExprOpt> ::= <Expr> | ε>

ForUpdate> ::= <AssignExpr> | ε>

"," <ReturnStmt> ::= "return" <ExprOpt>

"," <ExprStmt> ::= <Expr>

ExprList> ::= <Expr> <ExprTail> | ε>

ExprTail> ::= "," <Expr> <ExprTail> | ε>

<Expr> ::= <UnaryExpr> <RelOpTail> | <PointerRefExpr>

RelOpTail> ::= <RelOp> <UnaryExpr> <RelOpTail> | ε>

"!=" | "==" | "<=" | ">=" | "<" | ">" ::= <RelOp>

**UnaryExpr> ::= <UnaryOp> <UnaryExpr> | <PostIncrement> | <PreIncrement> | >
<<SimpleExpr**

PointerRefExpr> ::= "&" IDENTIFIER | "*" IDENTIFIER>

"--" PostIncrement> ::= IDENTIFIER "++" | IDENTIFIER>

PreIncrement> ::= "++" IDENTIFIER | "--" IDENTIFIER>

"!" | "-" ::= <UnaryOp>

<SimpleExpr> ::= <Term> <AddOpTail>

AddOpTail> ::= <AddOp> <Term> <AddOpTail> | ε>

"|" | "^" | "|" | "-" | "+" ::= <AddOp>

<Term> ::= <Factor> <MulOpTail>

MulOpTail> ::= <MulOp> <Factor> <MulOpTail> | ε>

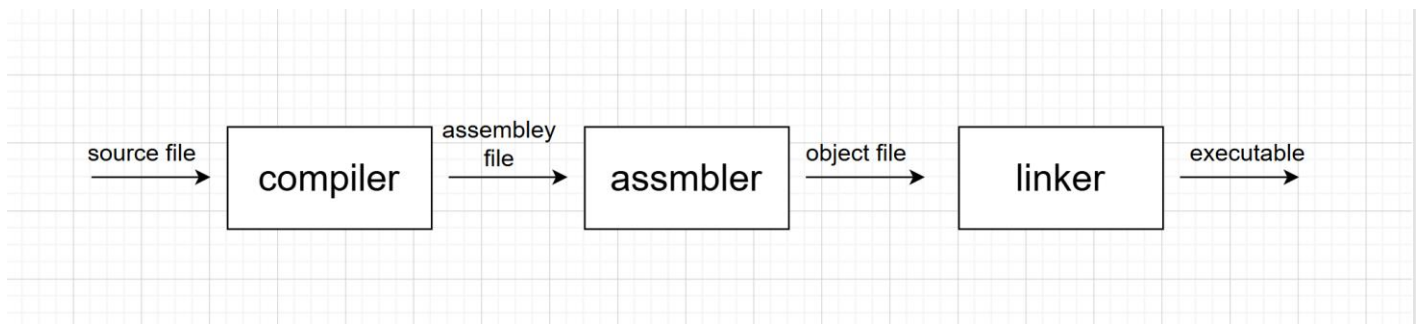
"&&" | "&" | "/" | "*" ::= <MulOp>

**Factor> ::= "(" <Expr> ")" | IDENTIFIER | INTEGER_LITERAL | FLOAT_LITERAL | >
"CHAR_LITERAL | IDENTIFIER "(" <ExprList> ")**

רקע תאורטי בתחום הפרויקט

קומפיילר/מהדר – הגדרה כללית:

בעולם המחשבים קומפיילר (או בשמו העברי: מהדר) הוא רכיב תוכנתי שתפקידו הוא להמיר בין שתי שפות תכנה, המהדר הקלאסי ימיר בין שפה עילית לשפת מכונה



ניקח לדוגמא את הקוד הבא הנכתב בשפה העילית "C" אשר נרצה להריץ:

```

1  int main()
2  {
3      int x = 10;
4      int y = 20;
5
6      x = x+y;
7
8      return 0;
9  }
```

שפות עליות כמו "C" ורבות אחרות עוצבו בצורה הדומה לשפה האנגלית המדוברת ושיהיה למתכנת נוח וקל להשתמש בהן אך המעבד לא בנוי בצורה שהוא יכול להבין בקלות את הוראות אלו. תפקידו של המעבד הוא להמיר את השפה העילית לשפת סף הבנויה בצורה קרובה לארכיטקטורת המעבד שאיתו אנחנו עובדים.

לצורך ההדגמה המהדר יצור לנו את קובץ ה assembly הבא:

```

1      .file      "dugma.c"
2      .def      __main;      .scl      2;      .type      32; .endif
3      .text
4      .globl    _main
5      .def      _main;      .scl      2;      .type      32; .endif
6      _main:
7      LFB0:
8          .cfi_startproc
9          pushl   %ebp
10         .cfi_def_cfa_offset 8
11         .cfi_offset 5, -8
12         movl    %esp, %ebp
13         .cfi_def_cfa_register 5
14         andl    $-16, %esp
15         subl    $16, %esp
16         call    __main
17         movl    $10, 12(%esp)
18         movl    $20, 8(%esp)
19         movl    8(%esp), %eax
20         addl    %eax, 12(%esp)
21         movl    $0, %eax
22         leave
23         .cfi_restore 5
24         .cfi_def_cfa 4, 4
25         ret
26         .cfi_endproc
27      LFE0:
28         .ident   "GCC: (MinGW.org GCC-6.3.0-1) 6.3.0"

```

אז נשתמש בעוד מהדר הממיר קוד סף אל שפת מכונה

(דוגמאות למהדירים כאלו: NASM לארכיטקטורת 64 x86 וTSEM לארכיטקטורת x86)

[illegible]

סוגים שונים של מהדרים:

Source compiler – ממיר את הקוד הראשוני (בשפה העילית) היישר לשפת מכונה לדוגמא -רוב המעבדים של השפות C CPP משתמשים בצורת הידור זאת

Intermediate Compiler - ממיר את הקוד הראשוני לשפת ביניים, אז שפת הביניים עוברת מהדר נוסף הממיר אותה לשפת מכונה. לדוגמא – המהדר של #C ממיר את קוד המקור לשפת CIL ומשם עובר עוד מהדר הממיר את הקוד לשפת מכונה

Transpiler – ממיר קוד משפה עילית אחת לשנייה לדוגמא – קוד בTypeScript יומר תחילה לJS בשלב הראשוני בהידורו (וזה נעשה בעזרת Transpiler)

סוגים שונים של שפות תכנות:

שפות תכנות פרוצדורליות - שפות כמו C ופסקל, בשפות אלו, התוכנית מורכבת מרשימת פקודות המבוצעות בסדר שלב אחרי שלב.

שפות תכנות מונחות עצמים -שפות כמו #C וJAVA. בשפות אלו, התוכנית בנויה ממודולים הנקראים אובייקטים, שכל אחד מהם כולל נתונים (תכונות) ופעולות (מתודולוגיות).

שפות תכנות פונקציונליות – שפות כמו Haskell וErlang. בשפות אלו, הפונקציה היא יחידת התכנות העיקרית. התכנות נעשה דרך קריאה לפונקציות, ולא על ידי שינוי מצב של משתנים.

שלבי הקומפילציה:

תהליך הקומפילציה מתחלק לשני שלבים מרכזיים,

ה-front end וה-beck end.

Front end:

שלב זה מתעסק בניתוח הקוד והמרתו לתצורת ביניים (IR) שאיתה השלב השני (beck end) יכול לעבוד. הוא עוסק בניתוח לוגי ובהבנה של המבנה והמשמעות של הקוד, וגם הוא מתחלק לכמה חלקים.

ניתוח לקסיקלי (Lexical Analysis)

השלב הראשון בfront end עוסק בחילוק הקוד הנתון לטוקנים, אשר כל טוקן הוא יחידה בעלת משמעות כמו מילים שמורות, משתנים, מספרים, אופרטורים וכדומה.

לדוגמה, פיסת הקוד `int x = 7;` תוכל להיות מתורגמת לטוקנים

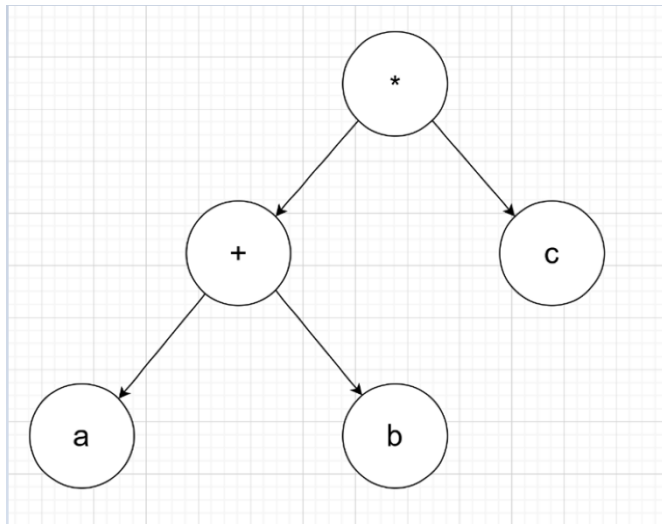
- `int` (מילה שמורה)
- `x` (מזהה)
- `=` (אופרטור)
- `7` (מספר)
- `;` (נקודה פסיק)

ניתוח תחבירי (Syntax Analysis)

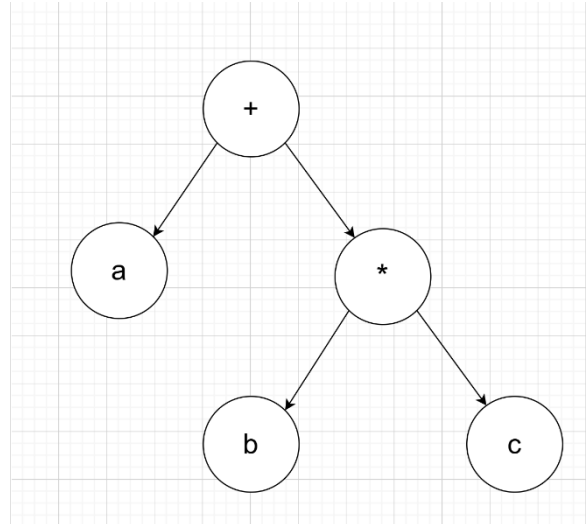
מטרה: השלב הבא הוא לבדוק אם רצף הטוקנים שנוצר בתהליך הקודם תואם לחוקי התחביר של השפה (הדקדוק שלה). השלב הזה יוצר את ה-עץ תחבירי או ה-עץ סמנטי (abstract syntax tree - AST), שמייצג את מבנה הקוד לפי כללי השפה.

דוגמה: אם הקוד הוא $a + b * c$, האנליזר התחבירי ייצור עץ תחבירי שמייצג את ההצהרה על המשתנה וההקצאה.

דוגמא לעץ סמנטי לא תקין:



דוגמא לעץ סמנטי תקין:



במהלך השלב הזה, אם יש בעיות תחביריות בקוד (כגון חוסר בסוגריים או סדר לא נכון של רכיבי השפה), יופיעו הודעות שגיאה.

דקדוק חסר הקשר (Context free grammar)

דקדוק חסר הקשר הינו ענף במתמטיקה אשר עוזר לנו להבין את המבנה של השפה והוא השלד שלפיו נבנה את העץ

סימני הדקדוק מתחלקים לשתי קבוצות **סופיים** ולא **סופיים**

סופיים: סימנים אשר ניתן להרחיב אותם לפי כלל דקדוק

לא סופיים: סימנים אשר לא ניתנים להרחבה

כללי דקדוק: דקדוקים בכללי ודקדוקים חסרי הקשר בפרט מורכבים מאוסף של חוקים אשר מגדירים את הדקדוק

כלל דקדוק בדקדוק חסר הקשר נראה ככה:

<רצף סימנים (סופיים או לא)> → <סימן לא סופי>

ניקח לדוגמא את הדקדוק הבא:

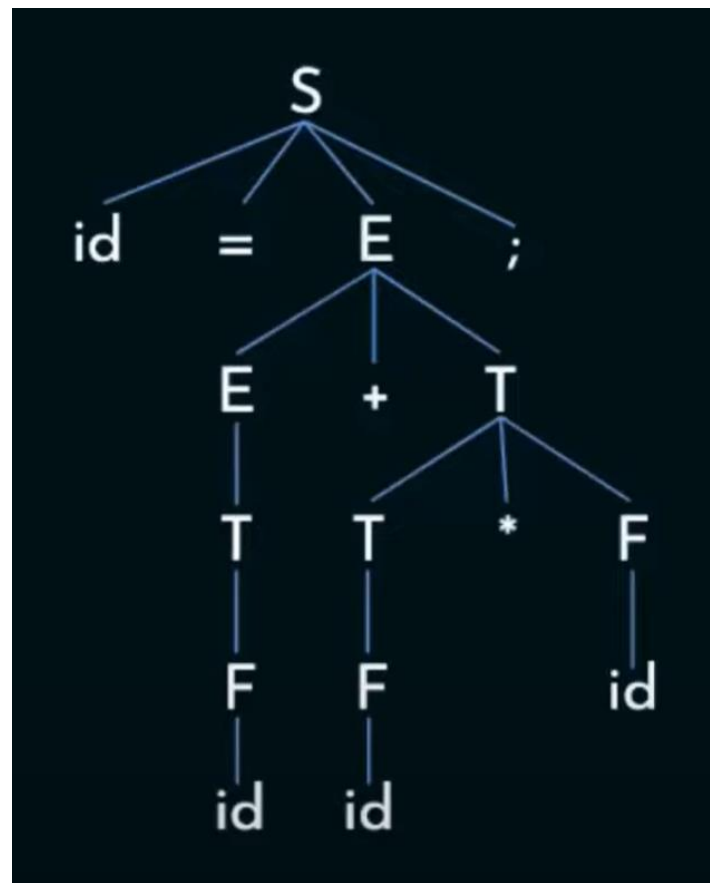
$$\begin{aligned} S &\rightarrow \text{id} = E ; \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow \text{id} \end{aligned}$$

הסימנים הסופיים בדקדוק: S,E,T,F

הסימנים הלא סופיים בדקדוק: id , = , ; , + , *

אם היינו מקבלים את הקלט הבא: **x = a + b * c;**

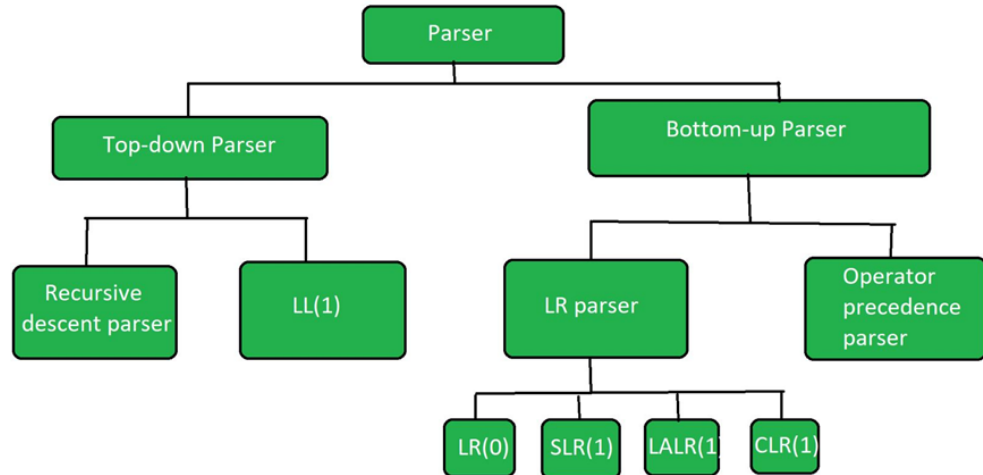
לפי הדקדוק הנתון יתקבל העץ:



סוגי מנתחים תחביריים

המנתח התחבירי הוא הרכיב המסובך ביותר במהדר. עקב זאת הוא בעל הכי הרבה תאוריה ומחקר בין כל הרכיבים

המנתחים התחביריים המתחלקים לשתי קבוצות עיקריות: **Top down** ו- **bottom up**



מלמעלה למטה (top-down)

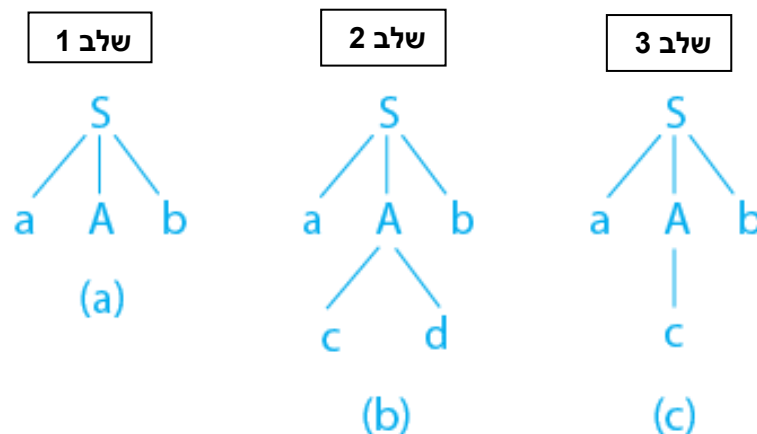
ניתוח תחבירי מלמעלה למטה מתחיל עם כלל התחבירי הראשון ומנסה לפתח אותו אל ה non terminals השונים בשפה

אם ניקח לדוגמא את ה grammar הבא:

$S \rightarrow aAb$

$S \rightarrow cd/c$

ניתוח top – down יוכל לפתח את העץ בצורה הבאה:



מנתחי bottom-up

מנתחי LR (left-right)

משפחה של מנתחים תחביריים מלמטה למעלה ומשתמשים באוטומט מחסנית pda

LR(0) – מנתח פשוט המשתמש הלא יודע להתמודד עם דקדוקים מסובכים ומורכבים מידי

LR(1) – עובד בדומה ל**LR(0)** אך מסתכל גם על האסימון הנוכחי וגם על האחד שאחריו וכך יכול להתמודד על דקדוקים מסובכים יותר

LALR(1) (**look-ahead, left-to-right**) – עובד בדומה ל **LR(1)** אך יוצר טבלה קטנה יותר וכך חוסך מקום זיכרון

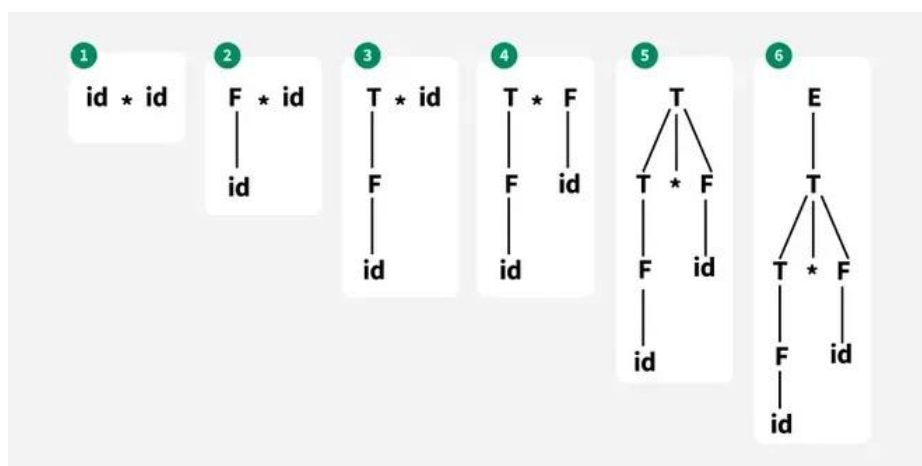
CLR(1) (**Canonical, left-to-right**) – המנטח הרובסטי יותר מכולם ויודע להתמודד עם כל הדקדוקים חסרי ההקשר שהם

מלמטה למעלה (bottom-up)

ניתוח תחבירי מלמטה למעלה ינסה לקחת את כל האסימונים ולצמצם אותם עד לכדי הגעה לכלל הראשון של הדקדוק

ינתח בצורה הבאה:

grammar הבא:



1. $E \rightarrow T$
2. $T \rightarrow T * F$
3. $T \rightarrow id$
4. $F \rightarrow T$
5. $F \rightarrow id$

Recursive Descent - מנתח המשתמש בפונקציה לכל כלל דקדוק

LL(k) Parsing – משתמש בטבלת ניתוח תחבירי ובמבט קדימה של k סמלים. יכול להיות ממומש עם מחסנית

ניתוח סמנטי (Semantic Analysis) –

שלב הניתוח הסמנטי הוא שלב שמטרתו לעבור על העץ הסמנטי ולוודא שאין בקוד שגיאות סמנטיות כגון, בדיקת טיפוסים, וידוא התאמה של פרמטרים לפונקציות ועוד....

-Back end

שלב זה עוסק בניתוח קוד הביניים שהתקבל בשלב הfront end. המטרה בשלב זה היא לא "הבנה" או "ניתוח" של קוד המקור אלא לעבור על קוד הביניים ולהעביר אותו בצורה אופטימלית לשפת הסף

אופטימיזציה –

שלב האופטימיזציה יעבור על קוד הביניים וייעל אותו בעזרת כמה שיטות, מחיקת שורות ללא משמעות, קיצור הקוד בדרכים שישאירו לו את אותה המשמעות וכו....

יצירת קוד הסף –

לקיחת קוד הביניים שעבר אופטימיזציה ולהמיר אותו לקוד סף

טיפול בשגיאות-

בכל שלבי ההידור נרצה לתפוס שגיאות.

לדוגמא, כתיבת שם משתנה לא תקין (שגיאה תחבירית) או פנייה למשתנה לא מוצהר (שגיאה סמנטית). טיפול בשגיאות הוא חלק קריטי בכל תהליך פיתוח תוכנה, ובפרט בקומפיילר. המטרה המרכזית היא להבטיח שהתוכנית שכתבת, או במקרה שלנו, הקוד שהקומפיילר מייצר, יתפקד בצורה נכונה, יציבה ובטוחה. נרצה להודיע למתכנת כשהקוד לא תקין ואיפה נמצאת השגיאה לנוחות ולפיתוח מהיר.

תיאור הבעיה האלגוריתמית

הבעיה האלגוריתמית טמונה במימוש ששת השלבי פיתוח הקומפיילר. כל שלב מציג בעיה אלגוריתמית שונה.

ניתוח לקסיקלי – חלוקת קוד המקור לטוקנים בעזרת אוטומט דטרמיניסטי סופי.

ניתוח תחבירי – שלב זה הוא מהעמוסים ביותר מבחינה אלגוריתמית בפיתוח המהדר, הוא כולל את יצירת העץ הסמנטי ועבודה אתו הכוללת מספר אלגוריתמים על עצים והכרעת מצבים בין בניית עצים שונים. שלב זה גם פועל על אוטומט דטרמיניסטי סופי.

ניתוח סמנטי – עבירה על העץ הסמנטי ובדיקה של התקינות הסמנטית שלו

יצירת קוד הביניים – עבירה על העץ הסמנטי ותרגומו להצגת הביניים.

אופטימיזציה – עבירה על קוד הביניים ובדיקה על איזה מקודות ניתן לעשות לו אופטימיזציה בדרכים שונות.

יצירת קוד הסף – תרגום קוד הביניים אל קוד הסף.

אסטרטגיה נבחרת לפתרון

האסטרטגיה – לפרק את הפרויקט הגדול לחלקים, כל חלק בעל פונקציונליות משל עצמו.

בסוף פיתוח כל חלק נשתמש בו כקופסא שחורה (black box) כך נוכל להפסיק להתעסק בחלק זה ולהתחיל לפתח את החלק השני בלי להתייחס לתאוריה של החלק הראשון, כך הפיתוח יהיה ממוקד יותר ומוציא קוד דינמי וקריא יותר.

להלן האסטרטגיה לפיתוח כל חלק:

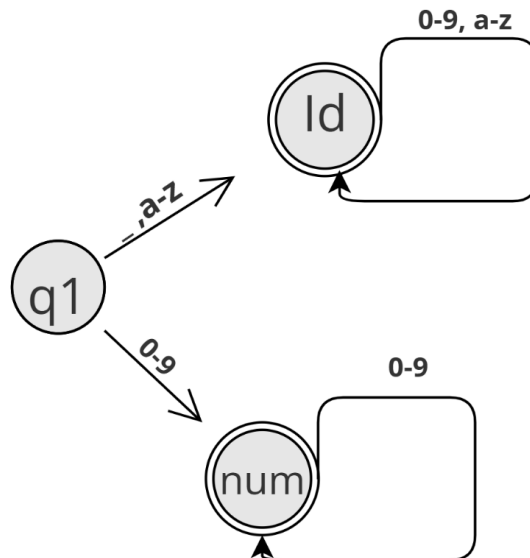
האסטרטגיה לפיתוח המנתח הסמנטי

בחרתי לפתח את המנתח הסמנטי בעזרת אוטומט סופי דטרמיניסטי (DFA). המנתח עובר תו תו על הקלט ומחלק אותו לאסימונים

ניתח לדוגמא שנראה לחלק את השפה לארבעה סוגים שונים של אסימונים:

מזהה , מספר , ואסימון לא תקין (לא שייך לשפה)

נבנה את האוטומט הבא:



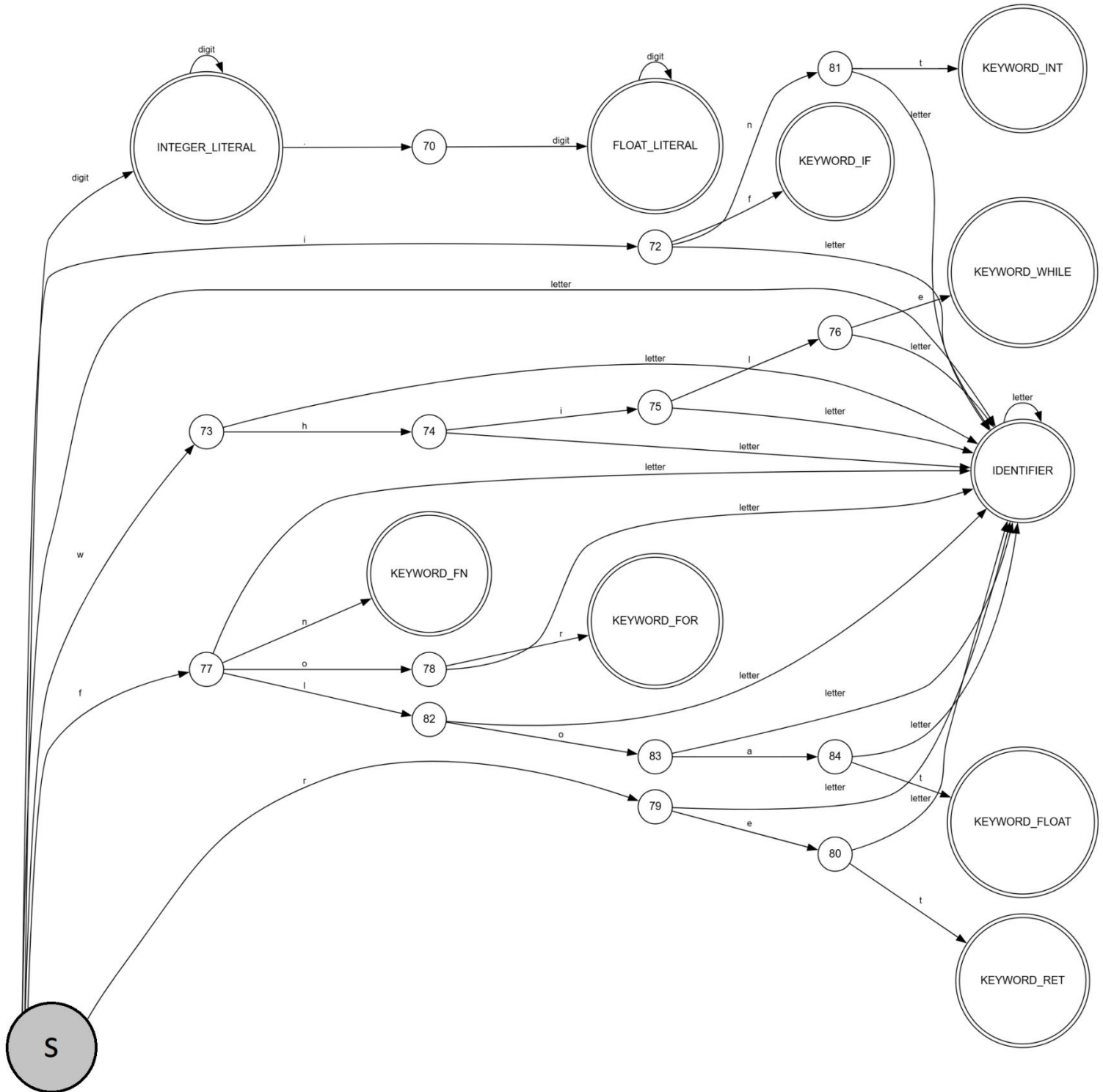
כפי שאנו רואים האוטומט למעלה יצליח לזהות מספרים ומזהים. כלומר: אם ניקח לדוגמא את המזהה sum נראה כי האוטומט יסיים בצומת Id וכך נדע שהוא מזהה. ציין כי עם קיבלנו קלט שהוא לא אחד מהקלטים האפשריים בצומת הנוכחית, נשלח לצומת מלכודת (Trap state) המציינת אסימון לא תקין

האוטומט של tec

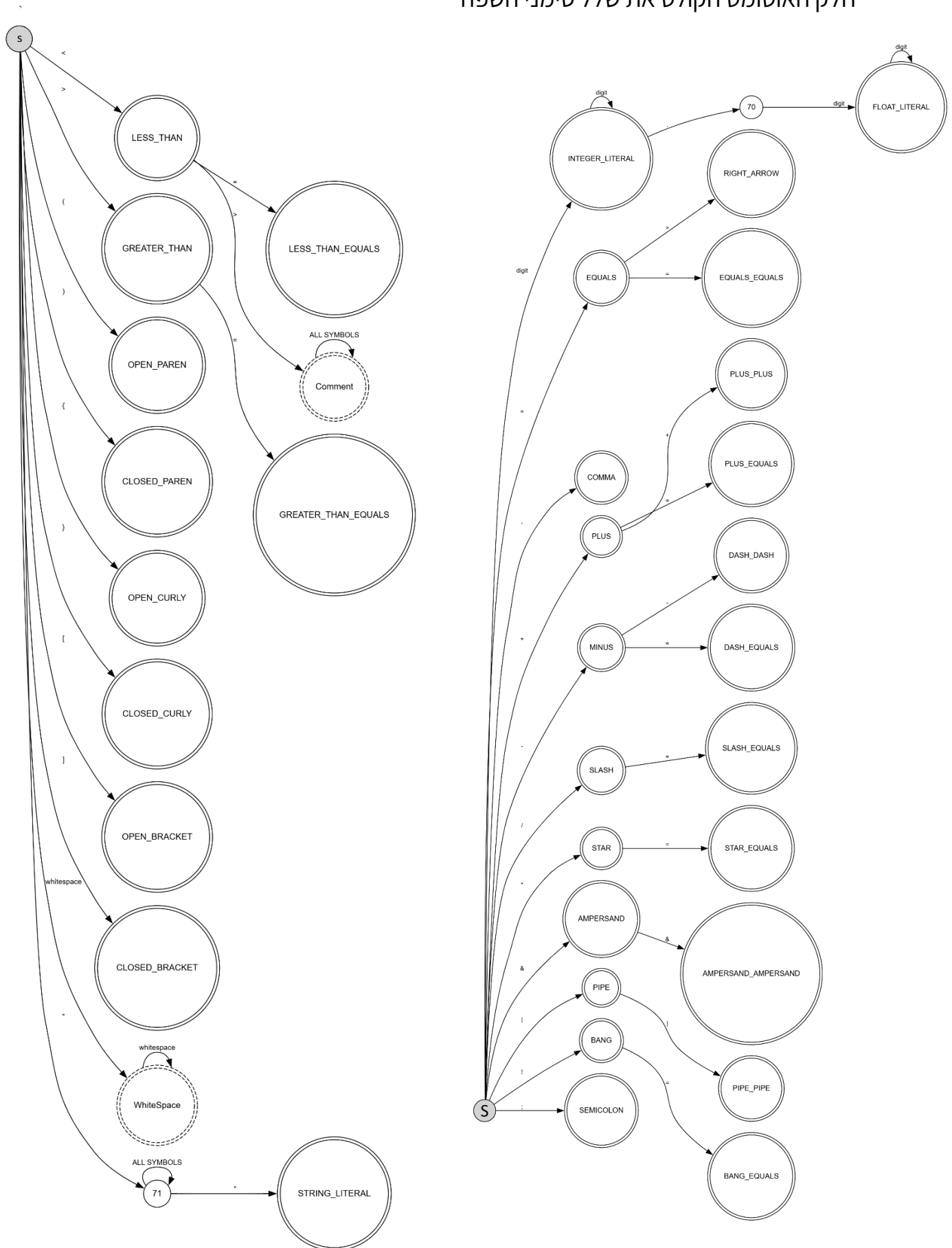
הסבר על התרשים:

- מצב התחלתי (מסומן ב-S) – מסמל את המצב הראשוני לסריקת אסימון לפני קבלת כל אות
- מצבי ביניים (מסומנים במספרים) – מסמלים מעבריי ביניים באוטומט. סיום הסריקה בהם תסתמן כשגיאה במהדר
- מצבי סיום (מסומנים בעיגול כפול) – מסמלים מצביי קבלה/סיום של האוטומט. סיום הסריקה בהם תסתמן כאסימון תקין
- מצבי דילוג (מסומנים בעיגול מקווקו) – מסמלים מצביי דילוג של האוטומט. סיום הסריקה בהם לא תחשב ככלום במהדר

חלק האוטומט הקולט מזהים (שמות משתנים ופונקציות),מחרוזות,מספרים שלמים,
מספרים עשרוניים ומילות מפתח:



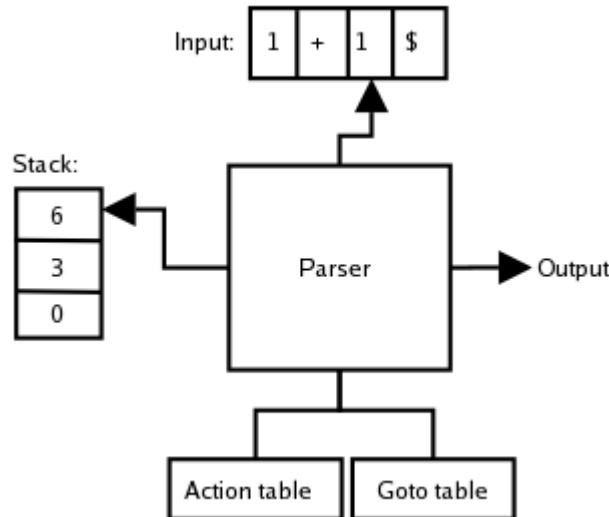
חלק האוטומט הקולט את שלל סימני השפה



האסטרטגיה לפיתוח המנתח התחבירי

בחרתי להשתמש ב LR parser לקומפיילר שלי

איך הוא בנוי?



המנתח התחבירי מסתמך על שתי טבלאות – **Goto, Action**

ומבנה נתונים – **מחסנית**.

הוא בנוי על מכונת מצבים גם כן

המחסנית – מאחסנת את סימני הדקדוק ומספרי המצבים

הטבלאות:

Goto – מסמנת לאיזה מצב לעבור על פי מצב הקלט והמחסנית

Action – מסמן איזה פעולה לבצע לפי המצב הנוכחי

הפעולות האפשריות

Shift – הכנס את האסימון הבא מהקלט אל המחסנית ביחד עם מספר המצב המצויין

Reduce – צמצם את הקלטים בראש המחסנית לפי חוק דקדוק מצוין

לפי מה נבנה את הטבלאות ?

ישנם כמה שיטות לבנות את הטבלאות אך אני בחרתי לבנות את הטבלה לפי אלגוריתם **IELR**

למה בחרתי באלגוריתם זה? כי הוא בונה את הטבלה יחסית מצומצמת ויכול לקבל הרבה סוגים שונים של חוקי דקדוק

ניקח לדוגמא את הCFG הבא:

(0)	$S' \rightarrow S$
(1)	$S \rightarrow C C$
(2)	$C \rightarrow c C$
(3)	$C \rightarrow d$

בעזרתו נבנה את טבלאות הACTION וגOTO

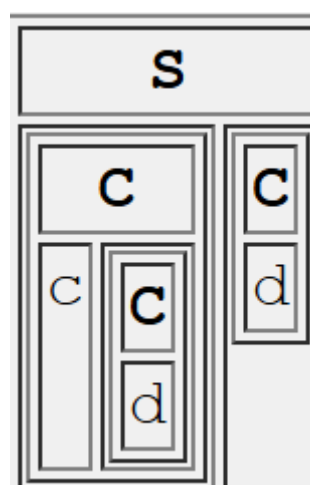
LR table						
State	ACTION			GOTO		
	c	d	\$	S'	S	C
0	s3	s4			1	2
1			acc			
2	s6	s7				5
3	s3	s4				8
4	r3	r3				
5			r1			
6	s6	s7				9
7			r3			
8	r2	r2				
9			r2			

ואם נרצה לנתח לדוגמא את הקלט: **c d d**

המעקב יראה ככה:

Step	Stack	Input	Action
1	0	c d d \$	s3
2	0 c 3	d d \$	s4
3	0 c 3 d 4	d \$	r3
4	0 c 3 C	d \$	8
5	0 c 3 C 8	d \$	r2
6	0 C	d \$	2
7	0 C 2	d \$	s7
8	0 C 2 d 7	\$	r3
9	0 C 2 C	\$	5
10	0 C 2 C 5	\$	r1
11	0 S	\$	1
12	0 S 1	\$	acc

והעץ הסופי יראה ככה:



האסטרטגיה לפיתוח המנתח הסמנטי

בשביל לממש שלב זה נצטרך:

לעבור על טבלת הסימנים ולראות ש-

- אין כפילויות.
- כל משתנה משומש לפי scope שלו

לעבור על העץ בצורה רקורסיבית, ולבדוק את התקינות של דברים שונים כגון:

האם בהכרזת וייסום משתנים סוגי המשתנים תואמים?

האם הפונקציה מחזירה את מה שנמצא בהכרזה שלה?

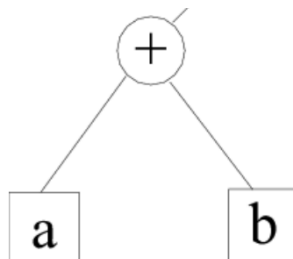
האם קריאה לפונקציה תואמת את הפרמטרים בהכרזה שלה?

ועוד...

האסטרטגיה לפיתוח שלב כתיבת קוד הסף

בשלב זה גם נעבור על העץ אך הפעם נסתכל על כל צומת שהיא לא בן בעץ ולפי סוגה נמיר לקוד סף

פעולות אריתמטיות



החלק בעץ:

יהפוך ל:

```

MOV RAX, [b]
ADD RAX, [c]
MOV [a], RAX
    
```

תנאים

נשתמש במבנה בסיסי של שפת סף לתנאים

הביטוי:

```
if (x == 5) {
    <גוף>
}
```

יהפוך ל:

```
mov     eax, DWORD PTR [x]
cmp     eax, 5
jne     .SOF
<התנאי גוף>
.SOF:
```

פונקציות:

יצירה:

האוגרים R9, R8, RCX, RDX, RSI, RDI – משמשים להעברת ששת הארגומנטים הראשונים של הפונקציה.

ניקח לדוגמא פונקציה.

```
int func1(int a, int b) {<גוף>}
```

נוכל ליצור אותה באמצעות יצירת label חדש וסיום בשימוש הפקודה ret

```
func1:
    <גוף הפונקציה>
    RET
```

כאשר a-b שמורים בRDI ו-RSI

קריאה לפונקציה:

```
CALL func1
```

לולאת for

```
for (int i = 0; i < 10; i++) {<גוף הלולאה>}
```

יתורגם ל:

```
MOV RAX, 0
MOV [i], RAX

for_start:
    MOV RAX, [i]
    CMP RAX, 10
    JGE for_end

    <גוף הלולאה>

    INC RAX
    MOV [i], RAX

    JMP for_start

for_end:
```

לולאות while

```
while (x > 0) {<גוף הלולאה>}
```

יתורגם ל:

```
MOV RAX, [x]

loop_start:
    CMP RAX, 0
    JLE loop_end

    <גוף הלולאה>

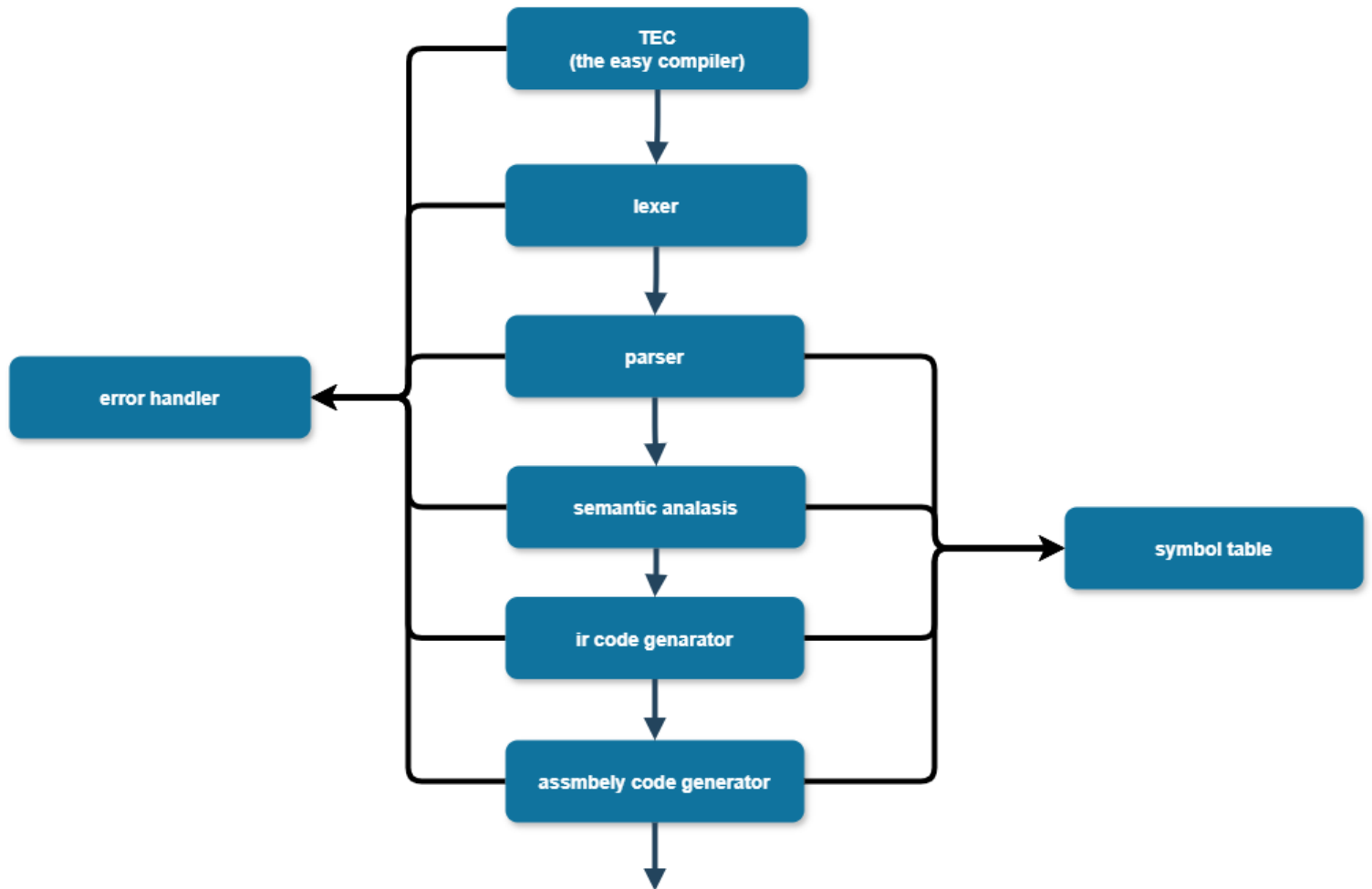
    MOV [x], RAX

    JMP loop_start
```

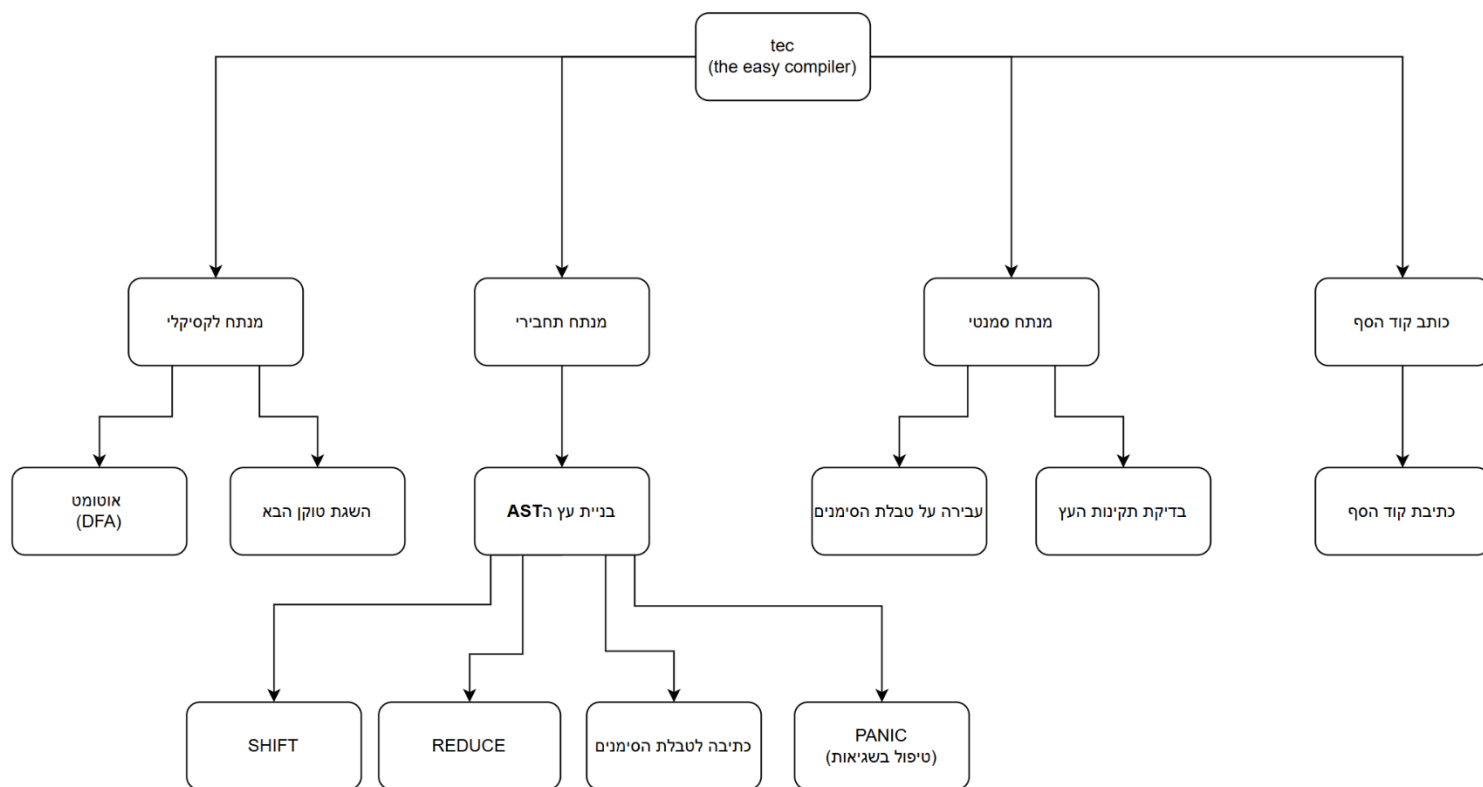
ארכיטקטורה של הפתרון המוצע בפורמט של

Top-Down Level Design

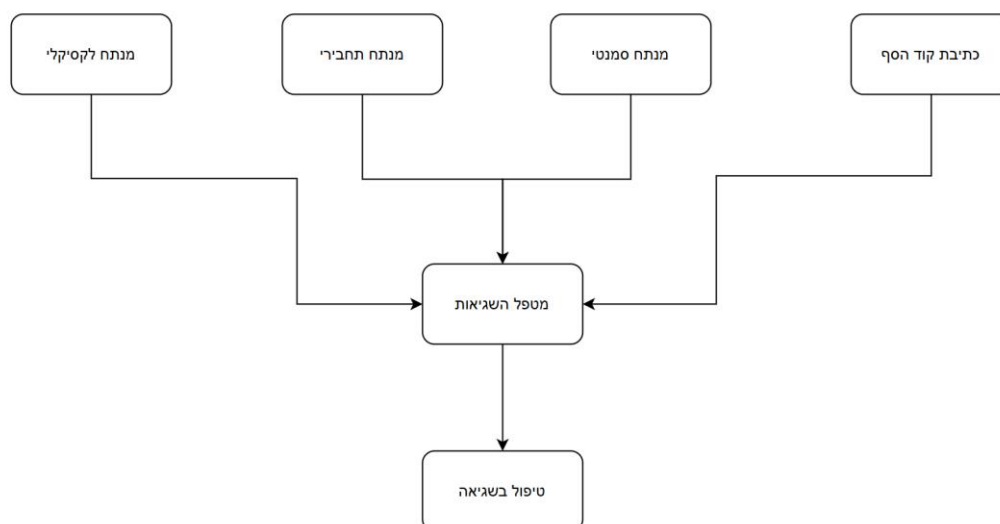
להלן תרשים זרימה המתאר את החלקים השונים ואת התקשורת ביניהם בפורמט השקפת על (High level):



תרשים Top-down מפורט



כל חלק אשר גזור ישירות מהקומפוננטה הראשית מקושר ישירות אל מטפל השגיאות



מבני נתונים

מטריצה

פיתוח מהדר כולל עבודה עם כלל מבני נתונים שונים ומגוונים. הרי אנחנו עובדים עם קלט אשר אנו רוצים לפרק ולעבד.

ובשביל לעשות זאת נשתמש **במבנה נתונים מתאים** לכל בעיה

כפי שראינו בפרקים הקודמים כתיבת המנתח הלקסיקלי כרוכה בעבודה עם **אוטומט סופי דטרמיניסטי (dfa)**. אך איך נמיר אותו לקוד.

ישנם שתי דכים מקובלות לעשות זאת.

מפת hash – אשר שומרת כמפתח את המצב ההתחלתי ובערך את הצמד של **תו המעבר והמצב הסופי**

מטריצה – אשר שומרת את אותם שלושה ערכים בפורמט שבו **השורות** מייצגות את המצב ההתחלתי **העמודות** את **תו המעבר והחיתוך ביניהן** הינו המצב הסופי

בחרתי להשתמש **במטריצה**

מטריצה הינו מבנה נתונים הבנוי כמערך דו ממדי של ערכים (בעל עמודות ושורות)

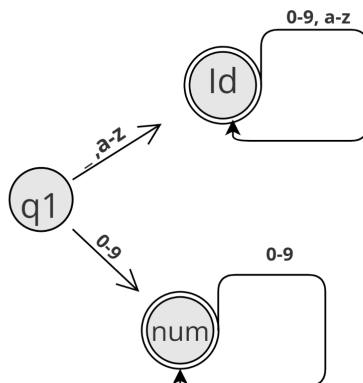
למה? – זהו מבנה נתונים אשר קל לעבוד איתו ולהכניס אילו נתונים.

שליפת **מצב הבא** בהינתן מצב נוכחי ותו מעבר הינו $O(1)$ לעומת **מפת hash** ששם היא $O(N)$ אך יש לשים לב שהמטריצה לא יותר טובה בצורה אבסולוטית מהמפה מכיוון שיעילות המקום של שהמטריצה תופסת בהתייחס לכמות המצבים שיש הינה $O(N^2)$ לעומת המפה אשר תופסת $O(N)$ מקום.

אך למרות זאת בחרתי במטריצה מכיוון שאעדיף להקריב **מקום** מאשר **זמן ריצה**

דוגמא להמרת אוטומט למטריצה

ניקח לדוגמא את האוטומט הבא:



נוכל להמירו למטריצה הזאת:

STATE	NUMBER	LETTER	-
q1	num	Id	Id
num	num	TRAP	TRAP
id	Id	Id	TRAP

קל מאוד לראות שאם לדוגמא מצבו הנוכחי של האוטומט הוא q1 והקלט הוא מספר נעבור למצב num

אנו נשתמש במטריצה בעוד מספר מקומות בקוד שבהם השימוש בא יותר טריוויאלי. כמו טבלאות (symbol table, goto table, action table)

מפת hash

Hash map הינו מבנה נתונים חשוב בתכנות אשר מאפשר לשמור צמדים של נתונים key,value בצורה כזאת שניתן אפשרות לשלף את ערך הvalue בעזרת המפתח ביעילות זמן ריצה $O(1)$

השתמשתי במבנה נתונים זה במקומות רבים בקוד שלי בעיקר לנוחות.

לדוגמא:

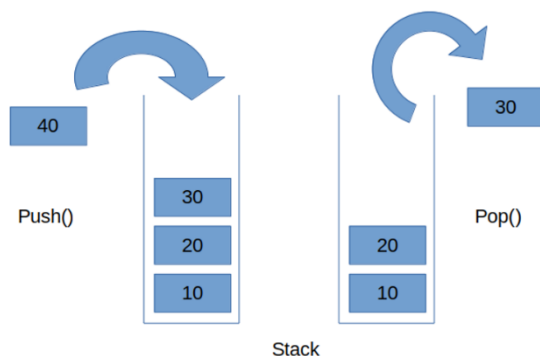
בטבלת הסימנים – שם משתנה/פונקציה (key) רשומה בטבלה (value)

באוטומט DFA - סימן בשפה (key) מיקומו במערך הסימנים (value)

מחסנית

מחסנית הינו מבנה נתונים העובד בשיטת Last in first out – LIFO (הראשון בפנים האחרון בחוץ)

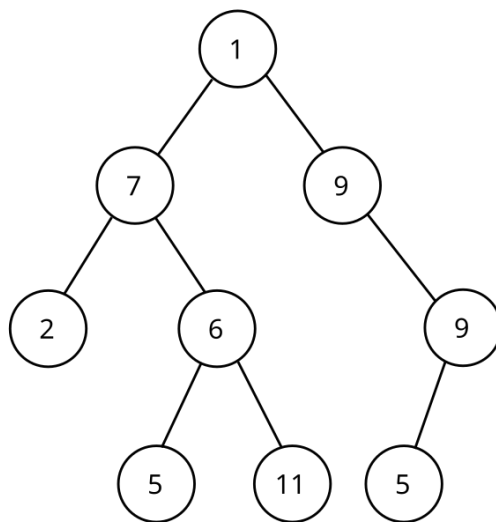
הוא מכיל שתי פעולות עיקריות push ו pop אשר משמשות להוצאה והכנסה של איברים מהמחסנית



נשתמש מחסנית כחלק ממימוש המנתח המילוני

עץ

עץ בינארי הוא מבנה נתונים הבנוי מצמתים, לכל צומת יש אפס "בנים" או יותר המחוברים אליו בקשת כך שלא יכולה להיות קשת המובילה מהבן בחזרה לאב (או לדרגה קדומה יותר)



דוגמא לעץ:

איך זה שימושי לנו?

החלק המרכזי של המנתח התחבירי הוא ליצור עץ המתאר את המבנה התחבירי של הקוד. הקוד בשפה בנוי בצורה היררכית ועץ זאת דרך מעולה לשמור את הנתונים בצורה בעלת בהמשמעות שאנחנו רוצים

מערך

מערך הינו מבנה הנתונים הנפוץ ביותר. ובכללי בתכנות, נשתמש בו בלי סוף

מדובר באיחסון הנתונים בצורה רציפה בזיכרון, דבר המאפשר התעסקות קלה ושליפת איברים ב $O(1)$

תיאור סביבת העבודה ושפות התכנות



לפרויקט זה, החלטתי להשתמש בשפת C++ מכיוון שהיא מאפשרת גם עבודת Low Level והתעסקות ישירה עם הזיכרון וגם שימוש בתכנות מונחה עצמים (OOP) אשר נותן לכל ישות בקוד להיות עצם משל עצמו בעל תכונות ופונקציות.

ככה אוכל לשמור על ביצועים גבוהים בעוד שימוש בתכנות מונחה עצמים וכלים מודרניים יותר



את הקוד עצמו כתבתי בתוכנת Visual Studio Code Version: 1.95.2



הידרתי אותו באמצעות GCC (gnu C compiler) אשר באופן אירוני בהכנת הפרויקט יצא לי לחקור עליו בעצמי וללמוד עליו

אלגוריתם ראשי

בפרק זה אתאר בפסיאודו קוד את האלגוריתם

אלגוריתם העל

1. אתחל את המנתח הלקסיקלי
2. אתחל את מדווח השגיאות
3. ייצר רשימת אסימונים מהקובץ הנתון
4. אם יש אסימון שגיאה דווח.
5. אתחל את המנתח המילוני
6. ייצר את העץ התחבירי
7. אתחל את המנתח סמנטי
8. עבירה על העץ הסמנטי
9. אתחל את מכולל קוד הסף
10. ייצר את קוד הסף

אלגוריתם המנתח הלקסיקלי

1. שמור את שם הקובץ
2. אתחל את טבלת האוטומט
3. עד שלא קיבלנו אסימון עם סוג סוף קובץ
3.1 הוסף לרשימת האסימונים את האסימון הבא

מציאת אסימון הבא

1. פתח את הקובץ
2. אתחל את המצב ההתחלתי
3. אם הגענו לסוף הקובץ
3.1 החזר אסימון סוף קובץ
4. עשה
4.1 קרא את התו הבא בקובץ
4.2 קבל את המצב הבא
4.3 אם המצב שקיבלנו לא תקין \ הגענו לסוף הקובץ
4.3.1 צא מהלולאה
4.4 אחרת
4.4.1 עדכן את הפוזיציה של האסימון ואת תכולתו
5. חזור ל 5.2
6. אם המצב שהגענו אליו הינו מצב דילוג
6.1 קרא את האסימון הבא
7. אחרת
7.1 אם זהו מצב סופי תקין של האוטומט
7.1.1 החזר אסימון חדש אם סוג מתאים למצב הסיום ועם התכולה שנצברה
7.2 אחרת
7.2.1 דווח שגיאה לקסיקלית למדווח השגיאות והחזר אסימון חדש עם סוג שגיאה

אלגוריתם המנתח התחבירי

1. אתחל את הטבלה
2. אתחל את חוקי הדקדוק
3. אתחל את קבוצת ה FOLLOW
4. אתחל את המחסנית
5. הוסף את מצב 1 לראש המחסנית
6. עשה
 - 6.1 הסתכל על המצב שבראש המחסנית
 - 6.2 קבל את הפעולה לפי האסימון הנוכחי והמצב שבראש המחסנית בטבלת ה ACTION
 - 6.3 אם פעולה הינה SHIFT
 - 6.3.1 בצע פעולת SHIFT
 - 6.4 אחרת אם הפעולה הינה REDUCE
 - 6.4.1 בצע פעולת REDUCE
 7. כל עוד הפעולה היא לא ACCEPT חזור ל 6
 8. הוצא מהמחסנית את האיבר הראשון (שורש העץ)
 9. החזר את השורש

Shift

1. דחוף את האסימון למחסנית
2. שנה את המצב הנוכחי לפי המספר בפעולה
3. עבור לאסימון הבא

Reduce

1. טעינת חוק הדקדוק התואם לפעולה
2. הוצאת מספר איברים לפי גודל חוק הדקדוק
3. ייצר צומת חדשה שכל ילדיה הם האיברים שהוצאנו והערך שלו הוא הסימן בצד שמאל של החוק
4. אם הצומת החדשה מסוג יצירת משתנה הוסף רשומה לטבלת הסימנים
5. דחוף את הצומת החדשה למחסנית
6. חשב את מצב החדש לפי סוג הצומת והצומת הקודמת עם טבלת ה GOTO
7. דחוף את המספר המתקבל

אלגוריתם המנתח הסמנטי

1. לכל צומת בעץ
 - 1.1 אם הצומת הינה מסוג השמה
 - 1.1.1 בדוק האם הסוגים בהשמה מתאימים
 - 1.1.2 בדוק בטבלת הסימנים האם הסימן קיים בScope הנוכחי
 - 1.2 אם הצומת מסוג קריאה לפונקציה
 - 1.2.1 בדוק כי הפונקציה קיימת
 - 1.2.2 בדוק כי הפרמטרים טועמים
 - 1.3 אם הצומת מסוג הכרזת פונקציה
 - 1.3.1 בדוק שהפונקציה לא קיימת
 - 1.4 אם הצומת מסוג return statement
בדוק שסוג ההחזרה תואם לפונקציה