

ANALIZOR AL GRADULUI DE COMENTARE PENTRU LINIILE DE COD INTEGRABIL ÎN SOLUȚII OPEN SOURCE

Radu Bucea - Manea - Țoniș

E-mail: radub_m@yahoo.com

Academia de Studii Economice, București

Rezumat: Se descrie proiectarea și implementarea unui analizor lexical, care generează documentația necesară optimizărilor clasice asupra codului C/C++. Analizorul conține un modul distinct de analiză a constantelor simbolice în vederea generării statisticilor de cuvinte cheie, pe baza cărora se realizează documentația primară a textului sursă. Datele furnizate în acest mod servesc analizelor calitative asupra codului, în măsură să justifice oportunitatea aplicării metodelor de optimizare. Rezultatele experimentale demonstrează că programele scanate cu analizorul lexical AGCC au un timp mai redus de compilare și generează mai puține erori.

Cuvinte cheie: compilator, analizor lexical, constantă simbolică, optimizarea codului.

Abstract: This paper describes the design and implementation of a lexical analyzer that automatically generates documentation to assist classic code optimizations. This scanner contains one new component, a distinct module for analyzing tokens and lexemes in order to produce the basic documentation of the source code, which is not commonly found in traditional lexical analyzers. The resulted documentation accumulates profile information and supplies this information to identify new optimization opportunities that are not visible to traditional methods. Experimental results show that the AGCC lexical analyzer significantly improves the performance of C Compilers in terms of execution time and error alerts.

Key Words: Compiler, lexical analyzer, token, code optimization.

1. Introducere în teoria compilatoarelor

Un compilator este un set de programe, care traduce un text în conformitate cu sintaxa unui limbaj evoluat de programare într-un fișier obiect scris în cod mașină. Acesta este transformat de un alt program denumit *linker* în fișier executabil.

Funcționalitățile unui compilator standard pot fi grupate în următoarele categorii de programe:

- **analizori lexicali** responsabili cu parcurgerea și transformarea primară a textului sursă;
- **pre-procesoare** responsabile cu expandarea codului original conform specificațiilor limbajului de macrocomenzi;
- **parsere** care analizează din punct de vedere sintactic codul intermediar și testează compatibilitatea cu standardul limbajului de programare;
- **analizorii semantici** verifică tipul variabilelor, compară definiția obiectelor cu referințele lor din program și emit avertizări;
- **generatorii de cod** transformă textul sursă validat în cod mașină după regulile limbajului de asamblare;
- **algoritmi de optimizare** care vizează programul la nivel de compilare, asamblare și execuție.

Operațiile clasice de optimizare globală a codului vizează conform [2]:

- 0.1. **optimizarea structurilor repetitive** prin eliminarea instrucțiunilor cu operanzi a căror valoare nu se modifică de la o iterație la alta;
- 0.1. **ștergerea instrucțiunilor** al căror rezultat nu este preluat de nici o altă secvență de cod;
- 0.1. **realocarea spațiului din memorie** alocat inițial variabilelor frecvent utilizate în zona de regiștri.

Calitatea compilatoarelor depinde de următoarele trăsături, conform [4]:

- **calitatea codului generat**, așa cum este definită în Standardul ISO/IEC 9126 - ca sumă a caracteristicilor de:
 - c₁. funcționalitate;
 - c₂. fiabilitate;
 - c₃. ergonomicitate;
 - c₄. eficiență;

- c₅. mentenabilitate;
- c₆. portabilitate.
- **detalierea mesajelor de eroare** în scopul identificării facile a erorilor de către dezvoltator;
- **viteza de compilare** depinde direct de calitatea și numărul proceselor de optimizare a codului, menționate ca opțiuni ale comenzii de compilare:
 - specificarea de biblioteci externe suplimentare;
 - generarea de cod destinat fluxurilor concurente de procesare;
 - identificarea liniilor de cod cu efecte imprevizibile;
 - generarea codului concomitent cu identificarea erorilor;
 - alocarea suplimentară de registre temporare codului compilat;
 - asigurarea unei portabilități maxime a codului;
 - generarea mesajelor de avertizare în cazul problemelor de compatibilitate.
- **integrarea cu instrumente de dezvoltare**, multe compilatoare fiind incluse în medii de dezvoltare ori colecții de biblioteci software, care diminuează semnificativ viteza de compilare;
- **asigurarea de suport** online atât pentru instalarea compilatorului, cât și pentru utilizarea sa în condiții optime;
- **conformitatea cu Standardul C++ ISO/IEC FDIS 14882:1998**; utilizarea de extensii adăugate de un producător și neconforme cu Standardul C++ conduce la reducerea portabilității codului și mărește dependența de un anumit tip de compilator; un compilator de calitate trebuie să accepte trăsături mai puțin uzuale ale Standardului C++ cum este cazul specializării parțiale a șabloanelor de funcții și clase;
- **prețul**, existând pe piață compilatoare liber download-abile și compilatoare contra cost, după cum rezultă din Tabelul 1.

Tabelul 1. Centralizator asupra compilatoarelor comerciale și Open source, după [4]

Compilatoare libere	Compilatoare contra cost
Apple C++	Borland C++
Bloodshed Dev-C++	CodeWarrior C++ (formerly Metrowerks)
Borland C++ v.5.5.	Comeau C++
Cygwin (GNU C++)	Edison Design Group C++ Front End
Digital Mars C++	Green Hills C++
MINGW - "Minimalist GNU for Windows"	HP C++ for Unix and HP C++ for OpenVMS
DJ Delorie's C++ development system for DOS/Windows (GNU C++)	IBM C++
GNU CC source	Intel C++ for Windows, Linux
Intel C++ for linux	Mentor Graphics/Microtec Research C++
The LLVM Compiler Infrastructure (based on GCC)	Microsoft C++
Microsoft Visual C++ 2008 Express edition	Paradigm C++
Sun Studio	The Portland Group C++
	SGI C++, optimizing compiler
	Sun C++
	WindRiver's Diab C++

Se recomandă compilatorul GNU GCC întrucât este un produs Open source, satisface majoritatea caracteristicilor de calitate și este compilatorul standard pentru sistemul de operare Linux.

2. Cerințe de calitate ale compilatoarelor moderne incluse în AGCC

Analizorul lexical AGCC a fost scris ca parte componentă a compilatorului GCC și moștenește de la acesta principalele caracteristici, conform cu [6]:

C.1. **analiza multi-pass a codului sursă**, presupune scindarea analizorului în mai multe secvențe de program, fiecare parcurgând o dată fișierul sursă și operând transformările necesare; principiul descompunerii compilatoarelor în subprograme după scopul funcțional a fost fundamentat la Universitatea Carnegie Mellon în cadrul proiectului Production Quality Compiler (PQC);

C.1.1. **împărțirea conceptuală a analizorului în Front-End și Back-End**, unde:

- **Front-End** construiește tabloul de cuvinte cheie și simboluri proprii limbajului de programare vizat, iar apoi parcurge fișierul în vederea construirii unui model intern al codului sursă; presupune următoarele etape:
 - e₁. **reconstrucția liniei**, în cazul AGCC, se realizează prin citirea caracter cu caracter și delimitarea viitoarei linii utilizând marcatori de sfârșit de linie;
 - e₂. **împărțirea codului sursă în constante simbolice** pe baza expresiilor regulate; această etapă mai poartă numele de scanare;
 - e₃. **pre-procesarea** în cadrul căreia AGCC își construiește un set propriu de directive care lansează în execuție procese paralele de analiză și transformare;
 - e₄. **construirea arborelui în vederea parsării** prin înlocuirea structurii liniare dictată de constantele simbolice cu una arborescentă, conformă cu sintaxa limbajului de programare.
- **Back-End** generează codul destinație împărțit în codul sursă parsat și fișierul de comentarii; generarea propriu-zisă este precedată de următoarele două faze:
 - f₁. **analiza** bazată pe reprezentarea intermediară a codului surprinde dependența dintre variabile; după [2], interesează colectarea datelor pentru aflarea numărului de:
 - apeluri individuale la funcțiile din program;
 - executări pentru fiecare bloc distinct de instrucțiuni;
 - validări ale condiției pentru fiecare ramură a unei structuri alternative și pentru fiecare caz din structurile condiționale de tip *CASE*.
 - f₂. **optimizarea generală a codului** vizează comasarea liniilor de cod, acolo unde este posibil, și eliminarea codului neoperațional.

Alte cerințe de calitate ale compilatoarelor moderne, după [5]:

- să se bazeze pe o **notație familiară** programatorului;
- **input minimal** necesar fluxului de prelucrări;
- **funcționalitate transparentă și predictibilă**;
- **dimensiune redusă și viteză de execuție mare**;
- **tratarea** corespunzătoare a **erorilor de parsare**;
- orientat pe **simplitate și eficiență**.

Analizorul AGCC satisface majoritatea cerințelor menționate anterior prin folosirea următoarelor tehnici, conform [1]:

- **algoritmi de sortare orientați pe cod**, economisind timpul necesar accesării tabelor; eficiența scanării depinde direct de lungimea codului sursă; următorul algoritm satisface acest criteriu urmărind frecvența de apariție a caracterelor "/" și "*":

```

if(pos[0]<=len){
    str_comp(i,0,pos[0]);
}
else if(pos[1]<=len){
    k=true;
    k=(pos[2]<=len)?false:true;
    if(k==false){
        str_comp(i,0,pos[1]);
        str_comp(i,pos[2]+2,len);
    }
    else {
        str_comp(i,0,pos[1]);
        return;
    }
}
else if(pos[2]<=len){
    k=false;
    str_comp(i,pos[2]+2,len);
}

```

- **citirea în blocuri de memorie** corespunzătoare sectoarelor de disc este mai eficientă decât introducerea fiecărui caracter în parte; din acest motiv AGCC permite încărcarea întregului text sursă în memorie pentru o parcurgere mai rapidă; se consideră că viteza de execuție a scanării crește în acest fel cu până la 30%;
- AGCC încurajează folosirea **instrucțiunilor de salt**, *goto*, întrucât nu își propune o analiză semantică asupra codului iar o prelucrare necesită o sigură parcurgere a textului sursă chiar dacă pot rezulta mai multe transformări.

3. Ciclul de dezvoltare al analizorului lexical AGCC

În teoria calculatoarelor, analizorii lexicali convertesc secvențe de caractere în constante simbolice. Un analizor este compus din scanner și generatorul de constante simbolice. Specificația limbajelor de programare prevede seturi de reguli de sintaxă, care precizează secvențele de caractere care pot forma constante simbolice. Constanta simbolică specifică tipul unui set de simboluri ori cuvinte cheie. Asocierea lor face obiectul funcțiilor de constante simbolice conținute de un analizor standard.

Din punct de vedere funcțional, AGCC înlătură toate comentariile din codul sursă original și le înlocuiește cu caracterul SPACE. Toate caracterele BACK-SLASH sunt șterse apoi, conducând la comasarea liniilor de cod detaliate în programul original pe mai multe rânduri pentru a obține mai multă claritate. AGCC identifică simbolurile și cuvintele cheie, după care stabilește tipul din care acestea fac parte. În această etapă de dezvoltare, analizorul ignoră simbolurile și cuvintele cheie, care nu au asociate o constantă simbolică. Scanarea secvențială a textului sursă se bazează pe tabelul de constante simbolice care în cazul AGCC se prezintă după modelul din Tabelul 2:

Tabelul 2. Cuvinte cheie și constante simbolice

Lexeme	Tokens
//	INLINE
/*	MULTI_BEG
*/	MULTI_END
;	SEMICOL
,	COMMA
)	R_BRACKET
"="	ASSIGN
{	L_CBRACKET
:	COLON
\	BACKSLASH
#	MACRO
struct	STRUCT
class	CLASS
void	VOID
bool	BOOL
int	INT
float	FLOAT
double	DOUBLE
char	CHAR

```

<Tipuri_date>
<Linie>
  <Nr>6</Nr>
  <Tip>class</Tip>
  <Den>Test</Den>
</Linie>
<Linie>
  <Nr>43</Nr>
  <Tip>int</Tip>
  <Den>testMe(int a)</Den>
</Linie>
<Linie>
  <Nr>51</Nr>
  <Tip>void</Tip>
  <Den>testMeToo(char c1)</Den>
</Linie>
<Linie>
  <Nr>57</Nr>
  <Tip>int</Tip>
  <Den>publicVar</Den>
</Linie>
</Tipuri_date>

```

Figura 1. Structura XML de constante simbolice

Căutarea unei constante simbolice încetează o dată cu întâlnirea de spații sau sfârșitului de linie. După identificarea constantelor simbolice, rezultatul este parsat în format XML și editat într-un fișier după modelul din figura 1. Pentru realizarea aplicației, se utilizează mediul integrat de dezvoltare Open source Eclipse, compilatorul GCC și colecția de biblioteci StandardTemplateLibrary care eficientizează scrierea de cod și menține viteza de execuție a programului. Codul sursă redactat conține metode denumite sugestiv și însoțite de comentarii clare. Executabilul obținut după compilare este referit direct. Deoarece fazele de analiză și optimizare pot necesita alocări suplimentare de memorie și timpi prelungiți de execuție, utilizatorul le omite în modul de lansare default (fără argumente suplimentare) al programului AGCC.

Se separă și liniile de comentarii însoțite de informații privind poziția în textul sursă, obținându-se cu ele un fișier distinct. Fișierul de comentarii este înglobat într-o structură portabilă de tip XML, care însoțește codul sursă sub forma unui manual de utilizare. În cazul produselor software de foarte bună calitate, dezvoltatorii includ funcția help care se bazează pe conținutul fișierului de comentarii, conținut pe care îl structurează în raport cu criteriile specifice aplicației.

Aceste informații permit:

- **construirea unei liste de variabile** la care se adaugă comentarii ce decurg din textul sursă privind tipul, domeniul, frecvența de utilizare în cod și legături cu alte variabile;
- **generarea comentariilor** pentru a marca semnificația secvențelor de prelucrare: calcul de sume, selecții dintr-o mulțime, construire de rapoarte și alte operații cu nivel de complexitate ce depinde strict de modul în care este proiectat produsul utilitar.

4. Testarea analizorului lexical AGCC

S-a impus o testare dinamică și analitică asupra AGCC, care presupune rularea unui program complet pentru executarea tuturor procedurilor și scenariilor de testare aferente diverselor stiluri în redactarea și comentarea codului C++. S-au aplicat aceleași proceduri de testare pentru a verifica comportamentul programului și pe alte sisteme de operare.

Se prezintă în Tabelul 3 output-ul analizorului lexical AGCC împreună cu textul sursă:

Tabelul 3. Testarea produsului AGCC

1. Comentarea integrală a unei linii de cod cu "//":	rezultă în urma parsării cu AGCC:
<pre>int main(void){ ... // Acesta este un comentariu ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 4: ... 5: return 0; 6: }</pre>
2. Comentarea parțială a unei linii de cod cu "/*":	
<pre>int main(void){ ... for(int i=0;i<10;i++){ // Acesta este un comentariu ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 3: for(int i=0;i<10;i++){ 4: ... 5: return 0; 6: }</pre>
3. Comentarea integrală a unei linii de cod cu "/*...*/":	
<pre>int main(void){ ... /* Acesta este un comentariu */ ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 4: ... 5: return 0; 6: }</pre>
4. Comentarea parțială a unei linii de cod cu "/*...*/":	
<pre>int main(void){ ... string city[3]={"Atena", "Bucuresti", "Sofia"}; /* Acesta este un comentariu */ ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 3: string city[3]={"Atena", "Bucuresti", "Sofia"}; 4: ... 5: return 0; 6: }</pre>
5. Comentarea mai multor linii cu "//":	
<pre>int main(void){ ... // Acesta este un comentariu // pe mai multe // linii. ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 6: ... 7: return 0; 8: }</pre>
6. Comentarea mai multor linii cu "/*...*/":	
<pre>int main(void){ ... /* Acesta este un comentariu pe mai multe linii. */ ... return 0; }</pre>	<pre>1: int main(void){ 2: ... 6: ... 7: return 0; 8: }</pre>
7. Comentarea pe mai multe linii cu "/*...*/" și "\":	
<pre>/* */ # /* */ def\ ne FO\ O 10\ 20 int main(void){ ... return 0; }</pre>	<pre>2: # define FOO 1020 7: int main(void){ 8: ... 9: return 0; 10: }</pre>

Testarea s-a efectuat asupra a 53 de linii de text sursă original rezultând 38 de linii de cod necomentat după o singură parcurgere a fișierelor sursă. Au fost eliminate 15 linii de comentariu, adică 28,3% din lungimea totală, fără să fie afectate structura și sintaxa proprie limbajului C++.

Există două excepții de care AGCC ține seama:

- e1. comentariile din cadrul directivelor de pre-compilare sunt ignorate;
- e2. asocierea de trigrafuri cu operatorul *backslash-new line* este evitată de AGCC prin convertirea acestora după eliminarea caracterului "\".

O practică bună în redactarea și depanarea codului în vederea executării optime a funcțiilor AGCC și compilării fără erori, presupune:

- verificarea **tipului parametrilor de intrare**, precum și **valorile întoarse de funcții**;
- **utilizarea șabloanelor** pentru eliminarea conversiilor nedorite între tipuri de date;
- **definirea de constante cu const** urmat de tipul de dată agreat;
- **declararea variabilelor** înainte de a le folosi și **eliberarea spațiului de memorie** alocat anterior imediat după ce au fost utilizate;
- **folosirea de paranteze** ori de câte ori logica calculului nu este foarte transparentă;
- **utilizarea assert** în tratarea excepțiilor pentru a preveni eventuale erori;
- folosirea unui **parametru de ieșire** pentru rezultat și **return** pentru eroare când o funcție trebuie să întoarcă o valoare calculată și o eroare în același timp;
- **verificarea valorilor de eroare** întoarse de funcții ale bibliotecilor de sistem pentru acele funcții care oferă acces la resursele sistemului, cum ar fi malloc() sau open().
- o **sistematizare a erorilor posibile**, bazată pe gradul de gravitate al acestora și pe influența lor diferită asupra rezultatelor finale ale programului.

Se recomandă evitarea:

- **conversiilor între tipuri de date diferite** și, mai ales, conversia pointerilor de un tip anume de date în void*;
- **definirii de tipuri derivate din pointeri** (ex: typedef char* Sir_de_caractere) și de constante preprocesare cu #define;
- introducerii mai multor **instrucțiuni pe aceeași linie** mai ales în cazul structurilor alternative și repetitive;
- **desfășurarea structurilor de control** pe mai mult de 3 niveluri;
- **invocarea exit()** din interiorul funcțiilor de bibliotecii;
- **folosirii goto, break sau continue** pentru ieșirea dintr-o structură și a structurilor repetitive post-condiționate de tip **do{...} while()**;
- **utilizării variabilelor globale** nesustenabile în cazul fluxurilor de prelucrare concurente și a limbajului de macrocomenzi.

5. Concluzii

Pe baza testelor desfășurate, se impun următoarele optimizări:

- o noua **schemă de alocare a memoriei** necesare pentru execuția AGCC;
- **documentarea comenzilor inline**;
- **rezolvarea** posibilelor **erori** sau **conflicte**;
- adăugarea în *Front-End* a facilităților pentru un **nou limbaj de programare**.

Se recomandă ca schimbările radicale să se producă în primele faze ale ciclului de dezvoltare astfel încât problemele care pot rezulta să fie remediate, iar adăugarea de patch-uri se va produce în urma unei analize cost-beneficiu. Se urmărește menținerea complementarității cu standardul C++ în vigoare în cazul versiunilor următoare.

Bibliografie

1. **MÖSSENBOCK, H.:** A Generator for Production Quality Compilers. Proc. of the Third International Workshop CC '90, Schwerin, FRG, October 22-24, Springer-Verlag, 1990, pp. 42-55. ISBN 3540536698, 9783540536697.
2. **POHUA, P., CHANG, SCOTT, A., MAHLKE, WEN-MEI, W., HWU:** Using Profile Information to Assist Classic Code Optimizations, Software - Practice and Experience, Vol. 21(12), 1991, pp. 1301-1321.
3. **STROUSTRUP, B.:** The Design and Evolution of C++, Addison-Wesley, 1994. ISBN 0201543303, 9780201543308.
4. * * *: <http://www.research.att.com/~bs/compilers.html/>
5. * * *: <http://gcc.gnu.org/onlinedocs/>
6. * * *: <http://en.wikipedia.org/wiki/>