

Latvijas Universitāte

Datorikas fakultāte

Ar regulārām izteiksmēm paplašinātu gramatiku dinamiska parsēšana.

Bakalaura darbs

Autors

Jūlija Pečerska

Vadītājs

Guntis Arnicāns

LU docents

Rīga, 2012

Anotācija

Anotācijas teksts latviešu valodā

Atslēgvārdi: Dinamiskās gramatikas, priekšprocesēšana, makro, regulārās izteiksmes, galīgi determinēti automāti, Python

Abstract

Abstract text in English

Keywords: Dynamic grammars, preprocessing, macros, regular expressions, determinate finite automata, Python

Saturs

1.	Ievads	4
2.	Darba pamatojums	6
3.	Transformācijas sistēma	9
3.1.	Transformācijas pieeja	9
3.1.1.	Priekšprocesori	10
3.1.2.	Dinamiskas gramatikas	11
3.1.3.	Parsētāja modifikācijas	11
3.2.	Parsētāju īpašības	11
3.3.	Makro sistēmas sintakse	12
3.4.	Transformācijas sistēmas apakšsistēmas	13
3.4.1.	Sakrišanu meklēšana	14
3.4.2.	Tokenu virkņu apstrāde	14
3.4.3.	Tipu sistēma	15
4.	Prototipa realizācija	18
4.1.	Atļautā makro sintakse	18
4.2.	Vispārīgā pieeja	18
4.3.	Makro konfliktu risināšana	19
4.3.1.	Divu makro konflikts vienā tvērumā.	19
4.3.2.	Divu makro konflikts dažādos tvērumos.	19
4.3.3.	Dažādu virkņu garumu konflikts	19
4.4.	Realizācijas pamatojums	20
4.5.	Kāpēc neder jau uzrakstītas bibliotēkas	20
4.6.	Lietotie algoritmi	21
4.6.1.	Regulāro izteiksmju pārveidošana nedeterminētā galīgā au- tomātā	21
4.6.2.	Determinizācija	23
4.6.3.	Minimizēšana	25
4.6.4.	Apvienošana	27
4.6.5.	Sakrišanu meklēšana	29
4.6.6.	Tvērumi	29
4.7.	Izņēmumi	30
4.7.1.	Transformācijas	30

4.7.2.	Produkcijas	30
4.7.3.	Tokenu klašu mantošana	30
4.7.4.	Regulārās izteiksmes daļu grupēšana	30
4.7.5.	Sapludinātā automāta minimizācija	30
4.8.	Optimizācijas iespējas	31
5.	Līdzīgu darbu apskats	32
6.	Rezultāti	33
6.1.	Prototipa testēšana	33
6.1.1.	Stresa testēšana	33
6.1.2.	Sistēmas testēšana	33
6.2.	Prototipa integrēšana Eq	35
7.	Secinājumi	37

Termini un apzīmējumi

Šeit būs aprakstīti termini, saīsinājumi un, ja būs nepieciešamība, apzīmējumi.

Regulārās izteiksmes

Atpakaļnorādes

Priekšprocesors

Tvērums Programmas tvērums ir programmas bloks, kurā definēti mainīgo nosaukumi vai citi identifikatori ir lietojami, un kurā to definīcijas ir spēkā. Programmas ietvaros tvērumus ievieš, piemēram, figūriekavas, C/C++ gadījumā. Tad mainīgie, kas tiek definēti vispārīgā programmas kontekstā (globālie mainīgie), var tikt pārdefinēti mazākajā kontekstā (piemēram, kaut kādas funkcijas vai klases robežās) un iegūst lielāku prioritāti. Tas nozīmē, ka ja tiek lietots šāds pārdefinēts mainīgais, tas tiek uzskatīts par lokālu un tiek lietots lokāli līdz specifiska konteksta beigām, nemainot globālā mainīgā vērtību.

Tvēruma piemērs:

```
int a = 0;
int b = 1;
int main() {
    int a = 2;
    a++;
    b += a;
}
```

Šajā piemērā *a* ir definēta gan globāli, gan lokāli. Kad tiek izpildīta rindiņa *a++*, lokāla mainīgā vērtība tiks samazināta uz 3, jo *a* ir pārdefinēts ar vērtību 2. Globālais *a* tā ar paliks ar vērtību 0. Un kad tiks izpildīta rindiņa *b += a*, *b* pieņems vērtību 4. Tvēruma iekšā tiks samainīta globālā mainīgā *b* vērtība, jo tas netika pārdefinēts.

Tipu izsecināšana

ϵ -pārejas

Nedeterminēts galīgs automāts

Determinēts galīgs automāts

1. Ievads

Основная цель этой работы - показать, что возможно сделать грамматику динамически парсируемой, контролируя её модификации. Данная работа представляет теоретическую базу для системы, которая позволит на лету расширять синтаксис языка. Представленная система является надстройкой над парсером, следовательно для её использования нет необходимости вносить значительные изменения в его работу.

Система разрабатывается в составе группы студентов университетов разных стран. Автор данной работы участвовал в разработке идеи в целом, однако фокусом была разработка подсистемы, которая позволит находить появления совпадения с шаблонами в коде программы.

В основе предложенной системы лежит принцип динамических - самомодифицирующихся - грамматик. Проблема с динамическими грамматиками заключается в том, что они могут много всего и контролировать их очень сложно. Однако данная работа выдвигает некоторые ограничения для контроля корректности изменений, вносящихся в грамматику. Для этого будет использоваться система вывода типов, при помощи которой можно будет следить за тем, чтобы переписывалось нужное и в нужное.

Данная система не позволит добавлять и удалять правила грамматики на усмотрение программиста. Однако она даёт удобную возможность дополнять существующие конструкции новыми при помощи языка макросов. Это существенно упростит работу с языком, так как чтобы добавить удобные конструкции не удет необходимости исправлять парсер языка.

Чтобы расширить возможности языка макросов, его система шаблонов будет расширена регулярными выражениями. Одна из особенностей описываемой системы заключается в том, что она работает не на тексте, а на токенах, за счёт чего и возможно устроить систему вывода типов. Далее список заматченных токенов из текста программы будет преобразовываться соответственно правилу трансформации, описанному в макросе.

В рамках данной работы будет описана модель предлагаемой системы. На данный момент система находится в стадии прототипирования, поэтому типовая система и система трансформаций будет только обозначена, а не рассмотрена полностью. Эта работа фокусируется на разработке эффективного прототипа для системы поиска совпадений в коде.

На момент написания работы есть готовый прототип системы поиска совпадений с шаблонами макросов. Основной принцип действия прототипа - создание одного детерминированного автомата для поиска совпадений по всем шаблонам по одному проходу по коду, по возможности за линейное время.

Далее организация работы будет следующей. Глава 2 рассказывает об актуальности проблемы, об общепринятом подходе к написанию парсеров языков и о возможностях динамических грамматик. Она также идентифицирует проблемы и сложности в использовании принципа самомодификации языка. Глава 3 рассказывает об основных принципах предлагаемой системы. Она даёт общий обзор языка модификаций и типовой системы, а также рассказывает о работе системы в рамках парсера языка. Глава 4 представляет прототип системы поиска совпадений с шаблонами, рассказывает об используемых алгоритмах и описывает причины выбора подходов. Там же описан подход к тестированию прототипа. В главе 5 описаны

результаты данной работы, глава 6 представляет выводы.

2. Darba pamatojums

Uz doto brīdi vairākumam mūsdienīgo programmēšanas valodu eksistē standarti. Pateicoties šim faktam, ir iespējams izstrādāt dažādus kompilatorus vienai un tai pašai valodai. Bet šiem standartiem ir jāmainās, lai valodas varētu attīstīties un paplašināt savas iespējas un lietojamību. Dažreiz šīs izmaiņas ir tik nopietnas, ka prasa ievērojamas parsētāja modifikācijas, lai tiktu atbalstītas, bet dažreiz tās ir tikai sintaktiskas, piem. sintaktiskā cukura¹ ieviešana.

Bieži izmaiņas ir valodas lietotāju iniciētas, jo tikai aktīvi lietojot valodu var izprast, kā tai trūkst, lai tā tiešām kļūtu par ērtu un populāru izstrādes rīku. Piedāvāt savas izmaiņas var dažādos veidos - informējot valodas izstrādātāju par kādas iespējas trūkumu, piedāvājot ielāpus kompilatoram, u.t.t.. Bet jebkādi izmaiņai ir nepieciešams atbalsts un lietotāju kvorums, kas vēlas šo izmaiņu lietot. To ir diezgan grūti dabūt ar šādiem izmaiņu pieprasījumiem.

Tomēr jebkādas, pat nebūtiskas, valodas sintakses izmaiņas parasti prasa kompilatora vai vairāku kompilatoru pārstrādi. Tas notiek tādēļ, ka programmēšanas valodas bieži iekļauj ļoti ierobežotas sintakses mainīšanas iespējas, vai arī neiekļauj šādu iespēju vispār.

Šīs darbs apskata iespēju lietot adaptīvo gramatiku principu, lai izvairīties no kompilatora pārstrādes nebūtisku sintakses izmaiņu gadījumā. Tas nozīmē, ka valodai jāsaturs konstrukcijas, kas parsēšanas laikā var modificēt un paplašināt pašas valodas sintaksi.

Viens no metodēm, kā varētu izpildīt pašmodificēšanas uzdevumu ir izveidot kross-kompilatoru, kas transformētu jauno sintaksi tā, lai standarta kompilators to varētu atpazīt. Bet šīs metodes problēma ir tas, ka lielākas daļas moderno valodu sintaksi ir neiespējams noparsēt lietojot automātiskos rīkus. Zemāk ir piedāvāti daži piemēri gadījumiem no populāras valodas C, kad automātiskā parsēšana ir neiespējama.

1. Valodā C lietotājs var nodefinēt patvaļīgu tipu lietojot konstrukciju `typedef`. Šāda veida iespēja padara neiespējamu šādas izteiksmes apstrādi $(x) + 5$, ja vien mēs neesam pārliecināti, kas ir x - tips vai mainīgais. Ja x ir tips, tad šī izteiksme pārveido izteiksmes $+ 5$ vērtību uz tipu x . Ja x ir mainīgais, tad šī izteiksme nozīmē vienkāršu mainīgā x un vērtības 5 saskaitīšanu.
2. Pieņemsim, ka ir iespēja paplašināt C valodas sintaksi ar infiksu operatoru `++` un pierakstīt konstanšu masīvus `[1, 2, 3]` veidā. Tad izteiksme `a ++ [1]` būtu nepārsējama, jo eksistē vismaz divi to interpretācijas veidi. Tas varētu tikt saprasts ka postfiksā operatora `++` pielietošana mainīgā `a` un tad `a` indeksēšana ar `[1]`. Vai arī tas varētu būt divu masīvu `a` un `[1]` konkatēnācija.

Dažreiz arī programmatūras koda dalīšana pa tokeniem ir atkarīga no šī koda konteksta, kas padara ne tikai parsēšanas procesu, bet arī leksēšanas procesu neautomatizējamu.

1. Valoda C++ ļauj lietotājam izveidot ligzdveida veidnes, piemēram, šādas
`template <typename foo, list <int>>`. Šajā gadījumā simboli `>>` aizver divas

¹ Sintaktiskais cukurs (*syntactic sugar*) ir speciālas konstrukcijas, kas tiek pievienotas lai to padarītu saprotamāku un lasāmāku cilvēkam. Šīs konstrukcijas nemaina valodas funkcionalitāti, bet gan atvieglo tās lietošanu. Izplatīts sintaktiskā cukura piemērs ir C valodas konstrukcija `a[i]`, kas patiesībā ir `*(a + i)`.

atvērtās grupas pēc kārtas. Lai tas tiešām būtu atpazīts, ka grupu aizvēršana, lekserim jāzin simbolu konteksts, jo parasti simbolu kombinācija >> nozīmē operāciju pārbīdei pa labi.

2. Gadījumā, ja lietotājam ir dota iespēja definēt savus operatorus, ieviešot operatoru, kas pārklāj eksistējošos, ir jāmaina leksēšanas likumus. Piemēram, ja lietotājs definē unāru operatoru $+-$, tad izteiksmei $+-5$ ir jābūt saprastai ka $(+-, 5)$, nevis ka $(+, -, 5)$.

Apskatītie piemēri dod iespēju secināt, ka automātisku parsētāju ģeneratoru lietošana var būt tik pat sarežģīta, ka parsētāja rakstīšana ar rokām. Izrādās, ka daudzām eksistējošām valodām parsētāji arī tiek rakstīti manuāli (piemēram C/C++/ObjectiveC kompilators GNU GCC [?]). Tas nozīmē, ka kross-kompilatoru visticamāk būs jāraksta manuāli, risinot eksistējošās gramatikas konfliktus, un oriģinālvalodas ievērojamu izmaiņu gadījumā būs jāpastrādā abi kompilatori, kas nozīmē divreiz vairāk darba.

Šis darbs apskata pieeju, kas lieto eksistējošo valodas parsētāju kā pamatu savam darbam un piedāvās likumu kopu, kas ļaus modificēt valodas gramatiku parsēšanas laikā. Taču patvaļīgas gramatikas izmaiņas var novest pie nekontrolējamās valodas evolūcijas. Tāpēc aprakstāmā pieejā tiek piedāvātas ierobežotas izmaiņu iespējas, kas tiks kontrolētas ar speciāli izveidotas tipu sistēmas palīdzību.

Sistēma ļaus ieviest jaunas konstrukcijas, konstruējot tās no jau eksistējošām valodas vienībām. Tā transformēs programmas gabalus attiecīgi pierakstītiem likumiem tā, lai valodas sākotnējā gramatika būtu tiem pielietojama. Šīs transformācijas korektumu nodrošinās tipu pārbaudes sistēma.

Ļoti līdzīgu uzdevumu, izņemot tikai transformāciju korektuma pārbaudi, pilda arī vispārējie priekšprocesori. Varbūt ir iespējams izveidot minēto sistēmu lietojot kādu no eksistējošām priekšprocesēšanas sistēmām, pievienojot tai kādas korektuma pārbaudes?

Jebkura priekšprocesora viens no galveniem mērķiem ir vienas elementu virknes aizvietošana ar citu. Virknes vienība var būt atšķirīga atkarībā no priekšprocesora, bet parasti šī vienība ir kādu vienas klases rakstzīmju kopa. Kļāšu daudzums parasti ir fiksēts (skaitlis, burts, tukšums, u.t.t.). Dažreiz zīmes piederība pie klases ir statiska, ka C priekšprocesorā, dažreiz ir dinamiska, ka, piemēram, \TeX , kur par atdalītāju var nedefinēt jebkādu specifisku simbolu. Tad apstrāde ir šo virkņu aizvietošana ar citām izveidotām virknēm.

Svarīgākā problēma šādai teksta apstrādes pieejai ir tas, ka tai neiespējams pārbaudīt korektumu. Apskatīsim sekojošu C makro piemēru:

```
#define foo(x, y) x y|
```

Pirmkārt, šādam makro nav iespējams statiski izsecināt rezultātu, jo kaut arī $foo(5, 6)$ tiks pārveidots par $5\ 6$, gan $foo(, 5)$, gan $foo(5,)$ tiks pārveidots par 5 . Otrkārt, tā kā komats ir makro daļa, nav iespējams kā pirmo makro argumentu padot virkni $5, 6$. To var izdarīt tikai ievietojot argumentu iekavās, tad $foo((5, 6), 7)$, kas tiks pārveidots par $(5, 6)\ 7$.

Gadījumā, ja ir nepieciešams saplacināt sarakstu, ir nepieciešams izveidot 2 makro, piemēram:

```
#define first(x, y) x
#define bar(x, y) first x y
```

Tomēr aprakstītais makro strādā tikai gadījumos, kad argumentiem ir pareizs tips. Piemēram, izteiksme `bar((5, 6), x)` tiks pārveidota par `5, x`. Bet izteiksme `bar(5, 6)` tiks pārveidota par `first 5 6`, kaut arī tai vajadzētu izraisīt kļūdu.

Var redzēt, ka vispārīgā gadījumā nav iespējams statiski izveidot nekādus secinājumus, tā kā makro rezultāts ir atkarīgs no argumentiem, kuriem tas ir pielietots. Bet patiesībā arī nekādus dinamiskus secinājumus nav iespējams izveidot, jo apstrādājot tekstu neeksistē korektuma kritēriji.

Tomēr pat neņemot vērā korektuma pierādījumu neiespējamību, makro sistēmai trūkst iespēju, lai izveidot jaunas valodas konstrukcijas. Piemēram, būtu dabiski reprezentēt skaitļa moduli ar pierakstu `|a|`. C priekšprocesors, savukārt, ļauj veidot tikai prefiksa formas funkciju makro un konstanšu makro. Jā arī kāda makro sistēma ļautu izveidot minēto pierakstu, parādītos problēmas gadījumos, kad vienam un tam pašam simbolam eksistē dažas nozīmes, piemēram, ar pierakstu `| (a | b) |`, kam jābūt pārveidotam uz `abs(a | b)`.

Apskatītie piemēri parāda, ka ne kross-kompilēšana, ne programmas teksta priekšprocesēšana vispārpieņemtā veidā neder vēlamā rezultāta sasniegšanai. Par šīs problēmas risinājumu varētu kļūt transformācijas sistēma, kas apstrādā nevis programmas tekstu, bet gan programmas tokenu un pseido-tokenu¹ virknes. Lai padarītu transformācijas sistēmu vairāk spēcīgu un ļautu atpazīt vispārīgākas virknes, tās makro šabloni tiks paplašināti ar regulāro izteiksmju elementiem, joprojām saglabājot iespēju kontrolēt transformāciju korektumu.

Šādā sistēma dos iespēju paplašināt valodu lietotāja līmenī, nevis kompilatora līmenī. Tas arī dos lielāku brīvību izmaiņu izstrādē, jo lietotājiem būs iespēja veidot makro bibliotēkas ar jaunām iespējām un izplatīt tās. Tā varēs būt pielāgota dažādām valodām, jo tā strādās ārpus paplašināmās valodas gramatikas.

Šī sistēma tiek izstrādāta zinātniskās grupas sastāvā, kurā piedalās cilvēki no Compiler Technology & Computer Architecture Group, University of Hertfordshire (Hertfordshire, England), Heriot-Watt University (Edinburgh, Scotland) un Moscow Institute of Physics and Technology (Dolgoprudny, Russia).. Tās ideja ir radusies valodas Eq² izstrādes darba gaitā un perspektīvā tiks integrēta ar šīs valodas kompilatoru.

¹ Pseido-tokens ir citu tokenu grupa, kas tiek aizvietota ar vienu objektu. Tas var tikt darīts, lai vienreiz noparsētu izteiksmi nevajadzētu apstrādāt vēlreiz. Tokenu aizvietošana ar pseido-tokeniem notiek gramatikas likumu reducēšanas brīdī. Kad, piemēram, tokenu virkne `id:a '+' id:b` tiek atpazīta ka derīga izteiksme gramatikas ietvaros, tā var tikt aizvietota ar pseido-tokenu `expr:a + b`.

² Atrodams tiešsaistē - <https://github.com/zayac/eq>

3. Transformācijas sistēma

Šī nodaļa piedāvā uzbūves principus sistēmai, kura ļaus lietotājam dinamiski paplašināt programmēšanas valodas iespējas ar makro valodas palīdzību. Šī makro valoda ļaus izveidot jaunas valodas konstrukcijas no jau eksistējošām vienībām. Sistēma ir projektēta ka virsbūve parsētājam, tā ļaus ieviest sintakses modifikācijas neieviešot nepieciešamību pēc parsētāja izmaiņām. Tā strādās paralēli ar parsētāju, analizējot kodu ar ierakstītiem makro un apstrādājot tokenu virknes.

Šīs sistēmas mērķis ir piedāvāt iespēju modificēt valodas sintaksi programmas rakstīšanas gaitā, nebojājot jau eksistējošo konstrukciju darbu. Sistēma ievieš pašmodificēšanos lietojot konstrukciju aizvietošanu, kas vienlaikus nodrošinās gramatikas modifikācijas un parsētāja nemainīgumu. Tajā pašā laikā sistēma būs stabila pret kļūdām dēļ tā, ka tā strādās tikai konkrētās gramatikas produkcijas ietvaros un tā, ka tā pārbaudīs tipus jaunizveidotām virknēm.

Sistēmas raksturīga īpašība ir tas, ka tā nav piesaistīta konkrētai programmēšanas valodai. Tā var tikt pielāgota un integrēta dažādu valodu parsētājos, kas atbilst dažiem nosacījumiem.

Tālāk nodaļa ir organizēta sekojoši. Apakšnodaļa 3.1. vispārīgi apraksta transformācijas pieeju. Apakšnodaļa 3.2. apraksta īpašības, kurām jāpiemīt parsētājam, lai tas varētu iekļaut transformācijas sistēmu. Apakšnodaļa 3.3. apraksta un pamato sistēmas makro sintaksi. Apakšnodaļa 3.4. detalizētāk apskata sistēmas sadalījumu uz trim apakšsistēmām un apraksta to sadarbību.

3.1. Transformācijas pieeja

Šī sistēma tiek projektēta ņemot vērā divu eksistējošo pieeju pieredzi. Pirmā no pieejām ir programmas koda priekšprocesēšana - koda makro ierakstu apstrāde pirms parsētāja darba sākšanos. Priekšprocesors parasti aizvieto kādas konstrukcijas ar citām noteikti definētam konstrukcijām. Otrā pieeja ir adaptīvās gramatikas - gramatikas, kas ļauj programmas kodam modificēt savu apstrādes gramatiku. Abas pieejas dod ļoti spēcīgus rīkus programmēšanas valodu izstrādē. Tomēr abām šīm pieejām ir savas problēmas un trūkumi, kurus šīs sistēmas projektēšanā mēģināja atrisināt.

Pirmā problēma, no kuras šī sistēmas izstrādē mēģināja izvairīties ir nekorekta simbolu ar divām nozīmēm apstrāde. Pieņemsim, ka mēs gribam funkciju $\text{abs}(x)$ apzīmēt ka $|x|$. Apskatīsim izteiksmi $| (a | b) + c |$, kurai vajadzētu tikt pārveidotai par $\text{abs}((a | b) + c)$. Gadījumā, ja transformācijas sistēma apstrādātu tekstu, tā nebūs spējīga pārveidot šādu konstrukciju. Tiks apstrādāti pirmie divi simboli, no $| (a |$, izveidojot nekorektu konstrukciju $\text{abs}((a | b) + c) |$. Šīs problēmas izvēlētais risinājums ir aprakstīts apakšnodaļā 3.1.1..

Otrā problēma ir dinamisku gramatiku nekontrolējamība. Dinamiskas gramatikas ir ļoti spēcīgs rīks, kas mūsdienās gandrīz netiek lietots. Tas tā ir tādēļ, ka dodot iespēju lietotājam patvaļīgi pievienot un dzēst gramatikas likumus, tiek zaudēta iespēja kontrolēt izmaiņu korektumu. Robežgadījums varētu būt tad, kad sākotnējā gramatika tiek pilnībā aizvietota ar citu gramatiku. Tad, kaut arī sākotnējā gramatika bija derīga parsēšanai ar eksistējošo algoritmu, jaunai gra-

matikai var piemīt īpašības, kas neļaus to apstrādāt (piemēram, kreisā rekursija LL parsētāju gadījumā). Šīs problēmas iespējama risinājums ir apskatīts apakšnodaļā 3.1.2..

Trešā problēma ir tas, ka dinamiskumu gramatikas līmenī nav iespējams ieviest jau eksistējošā valodas parsētājā bez ievērojamām izmaiņām. Parsētāja adaptācija, kas ļaus lietotājam modificēt gramatiku, ir sarežģīts darbs, kas dažreiz prasīs parsētāja arhitektūras izmaiņas. Iespēja pievienot valodai dinamiskumu nepārstrādājot parsētāju ir aprakstīta apakšnodaļā 3.1.3..

3.1.1. Priekšprocesori

Ir dažādas pieejas programmu pirmkoda priekšprocesēšanai. Visvairāk izplatītas no tām ir divas pieejas. Viena no pieejām ir sintaktiskā pieeja - sintaktiskie priekšprocesori tiek palaisti pēc parsētāja darbības un apstrādā sintaktiskos kokus, ko uzbūvē parsētājs. Dēļ aprakstāmās sistēmas īpašībām šajā darbā netiks apskatīti sintaktiskie priekšprocesori, jo līdz sistēmas darba izpildei parsētājs nevar uzbūvēt sintaktisko koku. Otra no pieejām ir leksiskā, leksiskie priekšprocesori tiek palaisti pirms pirmkoda parsēšanas un tiek nav zināšanu par apstrādājamas valodas sintaksi (piem. C/C++ priekšprocesors).

Leksiskie priekšprocesori pēc savām īpašībām ir tuvi aprakstāmai sistēmai. Ar makro valodu palīdzību tiem tiek uzdoti koda pārrakstīšanas likumi, un kods tiek pārveidots attiecīgi tām. Bet leksisko priekšprocesoru vislielākais trūkums ir tas, ka tie apstrādā tekstu pa simboliem neievērojot izteiksmju un konstrukciju struktūru. Apskatīsim jau minētu piemēru $\text{abs}((a \mid b) + c)$. Ar tādu makro sistēmu, kas neievēro koda struktūru, tātad neievēro to, ka patiesībā $(a \mid b) + c$ ir atomāra konstrukcija izteiksmē, šādu koda gabalu pareizi apstrādāt nevarēs¹.

Priekšprocesoru var iemācīt apstrādāt šāda veida konstrukcijas un atpazīt tos, ka atomārās izteiksmes. Bet tas nozīmēs, ka priekšprocesoram būs jāzina apstrādājamas valodas gramatika, kas neatbilst priekšprocesora lomai kompilēšanas procesā un nozīmē ka būs divreiz jāimplementē sintakses atpazīšana.

Lai izvairīties no šīs problēmas tika izvēlēts apstrādāt nevis programmas tekstu, bet gan programmas tokenu un pseido-tokenu virkni, ko daļēji jau apstrādāja parsētājs. Sistēma ir domāta ka virsbūve parsētājam un strādās paralēli ar to. Ielasītie makro tiks atpazīti tokenu virknē un pārveidoti citā tokenu virknē, kas aizvietos iepriekšējo. Tālāk aizvietotā virkne tiks apstrādāta ar sākotnējo valodas gramatiku. Konstrukciju $(a \mid b) + c$ parsētājs atpazīs ka pseido-tokenu `{expr}`, un sistēmas darbības laikā apstrādājamā plūsmā būs tieši pseido-tokens `expr`.

Otrā leksiskā tipa priekšprocesoru problēma ir tas, ka tie strādā ārpus programmas tvērumiem. Tas nozīmē, ka tvēruma sākuma tokens (piemēram, figūriekava, C/C++, Java un citu valodu gadījumā) tiek uzskatīts par parastu tekstu un var tikt pārrakstīts. Loģiskāk būtu, ja konkrētā tvērumā definēti makro tiktu mantoti līdzīgi ka mainīgie, kas nozīmē, ka šabloni, kas ir specifiski tvērumam, būtu ar lielāku prioritāti ka tie, kas definēti vispārīgākā tvērumā.

Sistēmas sakrišanas meklēšanas mehānisms tiks izstrādāts ņemot vērā programmas tvēru-

¹C/C++ priekšprocesors vispār neatļaus tādu konstrukciju izveidot, kaut arī šāds pieraksts ir diezgan loģisks no matemātiķu skatu punkta. C/C++ makro sistēma ļauj veidot tikai makro konstantes un prefiksa formas funkcijas.

ma maiņu. Tātad šabloni, kuri tiek ieviesti konkrētā tvērumā, strādās tikai tā ietvaros.

3.1.2. Dinamiskas gramatikas

Sistēma adaptē dinamisko gramatiku principu, ieviešot izmaiņu kontroli. Dinamiskas gramatikas vispārīgā gadījumā nekontrolē ieviestās izmaiņas, kas var sabojāt valodas parsējamību.

Transformācijas sistēmas makro nepievieno sistēmai jaunus gramatikas likumus tajā nozīmē, ka tie modificē gramatiku. Tie pievieno jaunu tokenu virknes pārrakstīšanas likumu, kas tiek izpildīts, kad tiek atrasta sakrišana ar makro specificētu šablonu. Tātad jaunās konstrukcijas tokenu virknē ir atpazītas un pārrakstītas uz konstrukcijām, kuras jau ir zināmas parsētājam. Makro sistēma nedos iespēju dzēst eksistējošos likumus no valodas gramatikas.

Lai nodrošinātu gramatikas likumu pievienošanas kontroli, tiek ieviestas dažas tipu specifikācijas. Katrs no makro attiecas uz kādu konkrētu gramatikas produkciju, un nevar tikt pārbaudīts citā parsēšanas brīdī. Katram makro pārrakstīšanas likumam arī ir specificēts tips, kas tiek lietots lai statistiski pārlicināties par to, ka transformācija ir korekta dotā tipa ietvaros. Tipu sistēma detalizētāk ir aprakstīta apakšnodaļā 3.4.3..

3.1.3. Parsētāja modifikācijas

Vēl viena svarīga dinamisku gramatiku īpašība ir tas, ka to lietošanai ir nepieciešams speciāli izstrādāts parsētājs, kas atļauj savu parsēšanas tabulu modifikācijas. Ir jāeksistē iespējai pievienot un dzēst attiecīgus gramatikas likumus, kā arī parsētājam jāprot pārbūvēt parsēšanas tabulas atbilstoši ieviestām izmaiņām.

Šī dinamisku gramatiku īpašība ļoti ierobežo iespēju lietot tos jau eksistējošo valodu paplašināšanai. Šāda veida papildināšana vairākumā gadījumu nozīmēs jauna parsētāja izveidošanu. Bet zinot, ka mūsdienīgām valodām parasti eksistē vairāki kompilatori, kurus izstrādā dažādi cilvēki, šādu izmaiņu ieviešana globālajā līmenī var kļūt neiespējama.

Aprakstāmā sistēma, savukārt, ir domāta ka palīgriks parsētājam. Parsētāju izmaiņas, kas būs nepieciešamas integrēšanai ar sistēmu ir minimālas. Tas dos iespēju lietot to jebkuram parsētājam, kas atbilst dažiem apstrādes nosacījumiem, kas ir apskatīti apakšnodaļā 3.2..

Sistēma dos iespēju papildināt valodas sintaksi bez nepieciešamības pilnībā pārstrādāt valodas parsētājus. Tā strādās ārpus valodas gramatikas. Pirms katras produkcijas apstrādes ar standartu valodas gramatiku likumiem tiek izsaukta transformācijas sistēma.

3.2. Parsētāju īpašības

Šajā darbā piedāvātā sistēma tiek izstrādāta uz LL(k) parsētāja bāzes. Lai parsētājs varētu tikt integrēts ar aprakstāmo sistēmu, tam jābūt implementētam ar rekursīvas nokāpšanas algoritmiem LL(k) vai LL(*). LL parsētājs tika izvēlēts tāpēc, ka LL ir viena no intuitīvi saprotamākām parsētāju rakstīšanas pieejām, kas ar lejupejošo procesu apstrādā programmatūras tekstu. LL parsētājiem nav nepieciešams atsevišķs darbs parsēšanas tabulas izveidošanā, tātad parsēšanas process ir vairāk saprotams cilvēkam un vienkāršāk realizējams.

Tā kā transformāciju sistēma tiek veidota ka paplašinājums parsētājam, parsētājam jāatbilst dažiem nosacījumiem, kas ļaus sistēmai darboties. Zemāk ir aprakstītas īpašības, kurām jāpiemīt parsētājam, lai tas varētu veiksmīgi sadarboties ar aprakstāmo sistēmu.

Tokenu virkne Parsētājam jāprot aplūkot tokenu virkni kā abpusēji saistītu sarakstu, lai eksistētu iespēja to apstaigāt abos virzienos. Tam arī jādod iespēju aizvietot kaut kādu tokenu virkni ar jaunu un ļaut uzsākt apstrādi no patvaļīgas vietas tokenu virknē.

Pseido-tokeni Parsētāji parasti pielieto (reducē) gramatikas likumus ielasot tokenus no ieejas virknes. Pseido-tokens, savukārt, konceptuāli ir atomārs ieejas plūsmas elements, bet īstenībā attēlo jau reducētu kaut kādu valodas gramatikas likumu. Viens no pseido-tokeniem, piemēram, ir tokens izteiksme - $\{expr\}$, kas var sastāvēt no daudziem dažādiem tokeniem (piem. $(a+b*c)+d$).

Vadīšanas funkcijas Pirmkārt, mēs prasam, lai katra gramatikas produkcija tiktu reprezentēta ar vadīšanas funkciju (*handle-function*). Ir svarīgi atzīmēt, ka šīm funkcijām būs blakus efekti, tāpēc to izsaušanas kārtība ir svarīga. Šīs funkcijas atkārtο gramatikas struktūru, tas ir ja gramatikas produkcija A ir atkarīga no produkcijas B, A-vadīšanas funkcija izsauks B-vadīšanas funkciju.

Sakrišanas funkcijas Katras vadīšanas funkcijas darbības sākumā tiek izsaukta tā sauktā sakrišanas funkcija (*match-function*). Sakrišanas funkcija ir transformācijas sistēmas saskarne ar signatūru (*Parser*, *Production*) \rightarrow *Parser*. Tā pārbauda, vai tā vieta tokenu virknē, uz kuru rāda parsētājs, ir derīga kaut kādai transformācijai dotās produkcijas ietvaros. Ja pārbaude ir veiksmīga, funkcija izpilda sakrītošās virknes substitūciju ar jaunu virkni un parsētāja stāvoklī uzliek norādi uz aizvietotās virknes sākumu. Gadījumā, ja pārbaude nav veiksmīga, funkcija nemaina parsētāja stāvokli, un parsētājs var turpināt darbu nemodificētas gramatikas ietvaros.

Ja izstrādājamās valodas parsētāja modelis atbilst aprakstītām īpašībām, tad tā var tikt veiksmīgi savienota ar aprakstāmo transformācijas sistēmu un ļaut programmētājam ieviest modifikācijas oriģinālās valodas sintaksē.

3.3. Makro sistēmas sintakse

Makro izteiksmes strādā stingri kaut kādas produkcijas ietvaros, tāpēc makro sintaksē tiek lietoti tipi, kas tiek apzīmēti ar produkciju nosaukumiem. Tipi tiks lietoti lai nodrošinātu pseido-tokenu virknes korektību sākotnējās gramatikas ietvaros pēc sintakses izmaiņu ieviešanas. Transformāciju sistēma sastāv no *match* makro likumiem un transformāciju funkcijām. Makro kreisā puse satur regulāro izteiksmi no tokeniem un pseido-tokeniem, kas tālāk tiek izmantota lai atrast virkni, kurai šī transformācija ir pielietojama. Makro labā pusē ir atrodamas funkcijas, kas izpilda transformē tās tokenu virknes, kas tiek akceptētas ar makro kreisās puses šablonu.

Apskatīsim *match* funkciju likumus, kas modificē apstrādājamās gramatikas produkcijas uzvedību. *Match* makro sintakses vispārīgo formu var redzēt attēlā 1..

`match [\prod1] v = regexp → [\prod2] f(v)`

1. att. *Match* makro sintakses vispārīgā forma

Šis apraksts ir uztverams sekojoši. Ja produkcijas *prod1* sākumā ir atrodama pseido-tokenu virkne, kas atbilst regulārai izteiksmei *regexp*, tad tai tiek piekārtots mainīgais ar vārdu *v*. Mainīgais *v* var tikt lietots makro labajā pusē kaut kādas funkcijas izpildē. Tātad ja tāda virkne *v* eksistē, tā tiek aizstāta ar pseido-tokenu virkni, ko atgriezīs *f(v)* un tālāk var tikt reducēta pēc gramatikas produkcijas *prod2* likumiem.

Regulārā izteiksme *regexp* ir vienkārša standarta regulārā izteiksme, kuras gramatika ir definēta attēlā 2..

```

regexp      → concat-regexp | regexp
concat-regexp → asterisk-regexp concat-regexp
asterisk-regexp → unary-regexp * | unary-regexp
unary-regexp  → pseudo-token | ( regexp )

```

2. att. Regulāro izteiksmju gramatika uz pseido-tokeniem

Pagaidām sistēmas prototipa izstrādē tiek lietota šāda minimāla sintakse, bet nepieciešamības gadījumā tā var tikt paplašināta.

Tagad mēs varam izveidot definētās makro sintakses korektu piemēru. Pieņemsim, ka ērtības dēļ programmētājs grib ieviest sekojošu notāciju absolūtās vērtības izrēķināšanai - `|{expr}|`. Sākotnējā valodas gramatikā eksistē absolūtās vērtības funkcija izskatā `abs({expr})`. Tad makro, kas parādīts figūrā 3. izdarītu šo substitūciju, ļaujot programmētājam lietot ērtāku funkcijas pierakstu.

```

match [{expr}] v = {|} {expr} {|}
      → [{expr}] {id:abs} {(} {expr} {)}

```

3. att. Makro piemērs #1

Vēl viens korektā makro piemērs: pieņemsim, ka funkcija *replace* ir definēta valodā *T* ar trim argumentiem, un darba gaitā tā jebkurā pseido-tokenu virknē aizvieto elementus, kas sakrīt ar otro argumentu, ar trešo funkcijas argumentu. Pieņemsim arī, ka mums ir nepieciešams izsaukt funkciju *bar* ar vienu argumentu, kas ir summa no funkcijas *foo* argumentiem. Šādā gadījumā makro, kas parādīts figūrā 4., izpildīs nepieciešamu darbību.

3.4. Transformācijas sistēmas apakšsistēmas

Šī nodaļa parāda sistēmas sadalīšanu uz trim neatkarīgam daļām. Pirmā no tām ir sakrīšanu meklēšanas daļa. Tā tokenu virknē atrod makro šablonu satikšanas reizes. Otrā ir atrasto


```
match [{expr}] v = {id:foo} {({} {expr} ( {,} {expr} ) * {})}
  → [{expr}] {id:bar} (replace v {,} {+})
```

4. att. Makro piemērs #2

virkņu apstrādes daļa. Tā pārveido sakrišanas mehānisma atrasto tokenu virkni atbilstoši tam, kas norādīts makro labajā daļā. Trešā ir tipu pārbaudīšanas sistēma, kas statistiski pārbauda, vai uzrakstītais makro var būt derīgs valodas gramatikas ietvaros. Šis sadalījums ir tikai konceptuāls, kas ir izveidots ērtības dēļ, lai varētu apskatīt sistēmu kā atsevišķu apakšsistēmu kombināciju.

Šablonu sakrišanu meklēšanas sistēma ir īsumā aprakstīta apakšnodaļā 3.4.1.. Sīkāk šīs apakšsistēmas īpašības un tās prototipa realizācija ir aprakstīta nodaļā 4..

Ir nepieciešams izveidot mehānismu, kas ļaus transformēt makro kreisās puses akceptētu pseido-tokenu virkni, izveidojot virkni, kas to aizvieto. Lai to izdarītu ir nepieciešama kaut kāds papildus rīks, par kuru ies runa apakšnodaļā 3.4.2..

Ir plānots, ka transformāciju sistēma varēs atpazīt nepareizi sastādītus makro šablonus lietojot tipu kontroli. Šīs pieejas bāzes principi ir aprakstīti apakšnodaļā 3.4.3.. Jāņem vērā tas, ka lai šī sistēma varētu tikt pielietota, izvēlētai transformāciju valodai jāpiemīt tipu secināšanas (*type inference*) īpašībai.

3.4.1. Sakrišanu meklēšana

FIXME: *Uzrakstīt!*

3.4.2. Tokenu virkņu apstrāde

Gadījumā, ja šablonu sistēma nesaturētu regulāro izteiksmju sintaksi (it īpaši *), būtu iespējams transformēt tokenu virknes ar pašu makro palīdzību. Atrastās tokenu virknes vienmēr būtu ar vienādu un determinētu garumu un saturu. Bet tā kā šabloni ļauj meklēt sakrišanas ar elementu sarakstiem, ir nepieciešams veids, kā apstrādāt jaunizveidotus un, iespējams, tukšus sarakstus.

Tātad lai varētu izpildīt atrastās tokenu virknes apstrādi un modificēšanu ir nepieciešams kaut kāds papildus rīks. Šis rīks varētu būt kaut kāda programmēšanas valoda. Izplatītākās programmēšanas valodu paradigmas mūsdienās ir imperatīvā vai funkcionālā paradigma. Katru no tiem varētu lietot atrasto virkņu apstrādei.

Šīm uzdevumam varētu lietot kādu no imperatīvam programmēšanas valodām, piemēram C, vienkārši izveidojot saskarni ar tās valodas kompilatoru. Bet vairākiem šādu valodu nav statistiskas tipu secināšanas iespējas. Statiska tipu secināšana C valodas gadījumā ir neiespējama rādītāju mainīgo dēļ, kuru tipus nevar izrēķināt parsēšanas laikā. Lai varētu ieviest statistiskās tipu izsecināšanas iespējas, vajadzēs ierobežot valodas iespējas, tātad modificēt eksistējošo kompilatoru vai kaut kā citādāk ierobežot pieejamo konstrukciju kopu.

Uzdevumam arī varētu lietot kādu no jau eksistējošām funkcionālām valodām ar tipu se-

cināšanas īpašību. Tad nebūs nepieciešamības veidot savu valodu pilnīgi no jauna. Bet tas nozīmēs, ka būs rūpīgi jāizpēta izvēlētais valodas sintaksi, kas var būt pārāk grūti.

Varētu lietot arī vienu no jau eksistējošām funkcionālām valodām ar tipu secināšanas īpašību, kas piemīt vairākam funkcionālo valodu. Vēl viena ērtā funkcionālo valodu īpašība ir tas, ka tās funkcijām nepiemīt blakusefekti, tātad to izpilde nevarēs samainīt eksistējošos datus. Valoda, kuras funkcijām ir blakusefekti, varētu sabojāt parsētāja darbu.

Šīs sistēmas implementācijā tika izvēlēts izveidot minimalistisku funkcionālu valodu, kura būs statistiski tipizējama. Tātad visiem šīs valodas mainīgajiem varēs izsecināt piederību pie tipa un pie kaut kāda virstipa, kas tiks lietots lai nodrošināt transformāciju korektumu. Funkcionālā pieeja nodrošina arī to, ka apstrādes funkcijām nepiemīt blakusefekti, kas varētu samainīt eksistējošos datus. Valodas mērķis ir ļaut izveidot atrasto tokenu permutācijas ar kādiem papildinājumiem nepieciešamības gadījumā.

Galvenais šīs valodas pielietojums ir dot iespēju apstaigāt pseido-tokenu virkni, kura tika atzīta par sakrītošu ar atbilstošu šablonu. Lai to darīt, tā dos iespēju lietot rekursiju un dažas iebūvētās funkcijas - saraksta pirmā elementa funkciju `head`, saraksta astes funkciju `tail` un objektu pāra izveidošanas funkciju `cons`. Funkcija `cons` funkcionālo valodu kontekstā strādā kā saraksta izveidošanas funkcija, jo saraksts `list(1, 2, 3)` tiek reprezentēta kā `cons(1, cons(2, cons(3, nil)))`, kur `nil` ir speciāls tukšs objekts. Valoda saturēs arī `if` konstrukciju, kas ļaus pārbaudīt dažādus nosacījumus.

Lai būtu iespēja apstādināt rekursiju, šī valoda arī ļaus izpildīt aritmētiskās operācijas ar veseliem skaitļiem. Tas dos iespēju izveidot skaitītājus un izveidot rekursijas izejas nosacījumus.

Tiek plānots, ka šī valoda arī ļaus izpildīt daļēju novērtējumu izteiksmēm, tur kur būs nepieciešams. Tas nozīmē, ka valodai jāsatur saskarne, kas ļaus piekļūt pie tokena vērtības. Šim mērķim ir domāta funkcija `value`, kas ir pielietojama pseido-tokeniem ar skaitlisku vērtību, piemēram, lai dabūt skaitli 5 no pseido-tokena `{int:5}`. Valoda arī ļaus izveidot jaunus tokenus ar izrēķinātu vērtību.

Funkcija `type`, savukārt, ļaus pārbaudīt tokenu tipu, kas var būt nepieciešams transformācijas procesā, piemēram, lai atpazīt kādu operatoru.

Lai būtu iespēja apstādināt rekursiju, šī valoda arī ļaus izpildīt aritmētiskās operācijas ar veseliem skaitļiem. Tas dos iespēju izveidot skaitītājus un izveidot rekursijas izejas nosacījumus.

3.4.3. Tipu sistēma

Kā bija redzams attēlā 1., katrā makro pusē ir atrodams produkcijas nosaukums, `[prod1]` un `[prod2]`. Tas tiek darīts tādēļ, lai kontrolētu, kad dotais makro ir pārbaudīts, un kāda tipa izejas virkni tas radīs. Abas šīs atzīmes ir rādītas tipu kontroles sistēmas dēļ.

Katrs atsevišķs makro strādā konkrētas gramatikas produkcijas ietvaros, `[prod1]` dotā makro gadījumā. Tas nodrošinās to, ka katrs no makro tiks izpildīts pareizajā vietā un visas konstrukcijas tiks apstrādātas.

Otrais tips, `[prod2]`, atzīmē to, ka pēc transformācijas procesa beigām mums jāsaņem tieši šādai produkcijai korektu izteiksmi. Tātad ir jāparbauda tas, ka funkcijas `f(v)` rezultāts

attiecībā uz atrasto tokenu virkni, ir atļauta ieejas virkne priekš produkcijas $prod2$.

Lai to paveikt ir nepieciešams izveidot pseido-tokenu regulāro izteiksmi produkcijai $prod2$. Tālāk ir nepieciešams izsecināt funkcijas f no virknes v rezultāta tipu.

Šajā darbā netiks apskatīts jautājums, kādā veidā tiks izveidota regulārā izteiksme priekš katras gramatikas produkcijas. To varētu izveidot programmētājs, vai, varbūt tā varētu tikt izveidota automātiski. Ir svarīgi pieminēt, ka pāreja no gramatikas likuma uz regulāro izteiksmi noved pie kādas informācijas zaudēšanas. Piemēram, nav iespējams uzkonstruēt precīzu regulāro izteiksmi valodai:

$$A := aAb \mid ab$$

Tomēr ir iespējams izveidot regulāro izteiksmi kas iekļaus sevī gramatikas aprakstīto valodu, piemēram, $a+b^+$. Makro lietotā transformācijas shēma tiks atzīta par pareizo, ja ir iespējams pierādīt, ka produkciju aprakstošā regulārā izteiksme atpazīst arī valodu, ko veido $f(v)$.

Ir viegli pamanīt, ka regulārās izteiksmes izveido dabisku tipu hierarhiju. Valoda, kura var tikt atpazīta ar regulāro izteiksmi r_1 , var tikt iekļauta citas regulārās izteiksmes r_2 atpazītās valodā apakškopas veidā. Piemēram, regulārās izteiksmes a^+ valoda ir atpazīstama arī ar regulāro izteiksmi a^* , bet a^* atpazīst vēl papildus tukšu simbolu virkni. Šādai tipu hierarhijai uz regulārām izteiksmēm eksistē arī super-tips, ko uzdod regulārā izteiksme $.^* - \top$. Ir acīmredzami, ka $\forall t_i \in R, t_i \sqsubseteq \top$, kur R ir visu regulāro izteiksmju kopa.

Ir svarīgi izveidot procedūru, kas ļaus izsecināt, vai $r_1 \sqsubseteq r_2$. Ir zināms, ka ir iespējams katrai regulārai izteiksmei uzbūvēt minimālu akceptējošu galīgu determinētu automātu¹. Šis automāts atpazīs precīzi to pašu valodu, ko atpazīst regulārā izteiksme. Tas nozīmē, ka no $r_1 \sqsubseteq r_2 \Rightarrow \text{semin}(det(r_1)) \sqsubseteq \text{min}(det(r_2))$. Diviem minimāliem automātiem A_1 un A_2 , $A_1 \sqsubseteq A_2$ nozīmē, ka eksistē kaut kāds attēlojums Ψ no A_1 stāvokļiem uz A_2 stāvokļiem, tāds, ka:

$$Start(A_1) \rightarrow Start(A_2) \in \Psi$$

$$\forall s \in States(A_1) \forall e \in Edges(s), \Psi(Transition(s, e)) = Transition(\Psi(s), e)$$

Šeit $States(x)$ apzīmē automāta x stāvokļu kopu, $Edges(s)$ apzīmē pseido-tokenu kopu, kas atzīmē no stāvokļa s izejošās šķautnes. $Transition(s, t)$, savukārt, apzīmē stāvokli, kas ir sasniedzams no s pārejot pa šķautni, kas atzīmēta ar pseido-tokenu t .

Otra svarīga īpašība, kas tiks lietota šajā tipu pārbaudīšanas sistēmā ir tas, ka transformāciju valoda ir statiski tipizējama un tā satur ļoti ierobežotu iebūvēto funkciju skaitu. Katrai no šīm iebūvētām funkcijām ir iespējams izveidot to aprakstošo regulāro izteiksmi. Piemēram, regulārā izteiksme funkcijai $head(x)$ var tikt izveidota ka visu to šķautņu kopa, kas iziet no x aprakstošā automāta sākuma stāvokļa. Kad šablona apstrāde tiek sākta, var arī izsecināt iespējamo virknes garumu intervālu un izvadīt brīdinājumus gadījumā, ja ir iespējama funkcijas $head(x)$ izsaukšana no tukšā saraksta.

¹Sk. nodaļu 4.6..

Šī nodaļa tikai vispārīgi apraksta tipu secināšanas sistēmas ideju. Uz doto brīdi tā atrodas izstrādes stadijā, bet var redzēt ka tās izveide ir pamatota. Sīkāka informācija par tipu sistēmu ir saņemama pie Eq kompilatora izstrādes komandas.

4. Prototipa realizācija

Šī darba ietvaros tika izstrādāts prototips sakrišanu meklēšanas sistēmai. Šī nodaļa apraksta prototipa īpašības, kā arī pieejas un algoritmus, kas tika lietoti tā realizācijā. Prototips vienkāršības un izstrādes ātruma dēļ tika rakstīts Python valodā, un ir viegli palaižams un atklūdojams uz jebkura datora ar pieejamu 2.7. Python versiju.

Prototips tika izstrādāts ar iedomu pēc iespējas samazināt sakrišanu meklēšanas laiku, jo transformācijas sistēmas izsaukumi notiks katras produkcijas apstrādē.

Šīs nodaļas organizācija ir sekojoša. Apakšnodaļa 4.1. īsi apraksta prototipam pieejamu sintaksi. Apakšnodaļa 4.2. apskata darbību virkni, ko prototips izpilda darba gaitā. Tālāk apakšnodaļa 4.3. apskata izvēlētās metodes šablonu konfliktu risināšanai. Apakšnodaļā 4.4. ir aprakstīts, kāpēc tika izvēlēta šāda realizācijas pieeja un apakšnodaļa 4.5. apskatīts, kāpēc dotajā gadījumā nav derīgi kādi gatavie risinājumi. Apakšnodaļa 4.6. stāsta par realizācijā lietotiem algoritmiem, bet apakšnodaļa 4.7., savukārt, apraksta problēmas ar kurām saskārās darba autors un izņēmumus, kas pagaidām netiek implementēti prototipā.

4.1. Atļautā makro sintakse

Kā jau bija pateikts, makro pieejamā regulāro izteiksmju sintakse ir minimāla. Tā atļauj lietot `*` lai identificēt tokenu virknes un `|` lai izvēlēties starp dažiem tokenu tipiem.

Prototips arī ļauj veidot regulārās izteiksmes ar specifisku tokenu vai pseido-tokenu vērtībām. Piemēram, regulārā izteiksme `{id:foo}` sagaidīs tieši identifikatoru `foo`, bet izteiksme `{id}` sagaidīs jebkuru identifikatoru. Tas ievieš dažādas problēmas, kas tiks aprakstītas zemāk, bet dod lielāku brīvību šablonu sistēmas lietotājam.

4.2. Vispārīgā pieeja

Prototips imitē darbu reālajā vidē, saņemot tokenus no ieejas plūsmas pa vienam no atsevišķas saskarnes. Tokenu plūsmas saskarne imitē leksera darbu. Kamēr prototips nav saņēmis nevienu regulāro izteiksmi, tas ignorē tokenu plūsmas apstrādes izsaukumus. Tiklīdz prototipam atnāk izsaukums apstrādāt tokenu regulāro izteiksmi, tas uzsāk regulārās izteiksmes parsēšanu. Parsēšanas procesā tiek izveidots galīgs automāts, kas akceptē regulārās izteiksmes uzdotās virknes. Katra jauna regulāra izteiksme izveido jaunu automātu.

Tad, kad atnāk tokenu apstrādes pieprasījums, sistēma izpilda pārejas starp automātu stāvokļiem, meklējot sakrišanas, un atceras tokenus, kurus jau ir nolasījusi. Sistēma atrod garāko virkni, kas atbilst kādam no šabloniem un tad atgriež tās identifikatoru un nolasīto tokenu virkni, lai turpmāk transformēšanas mehānisms varētu pārstrādāt to jaunajā virknē.

Pieņemot, ka transformēšanas sistēma ir izstrādāta, tālākā darba gaita būs sekojoša. Transformēšanas sistēma aizstāv ielasīto virkni ar citu, kas ir konstruēta pēc akceptētā makro šablona noteikumiem. Tad sakrišanas meklēšanas sistēmas darbs tiek uzsākts no jauna no aizvietotās virknes sākuma.

Sistēma turpina darbu aprakstītā gaitā līdz ko neviens no šabloniem vairs netiek akceptēts. Pēc sistēmas apstāšanās tiek iegūta jauna tokenu virkne, kas tika apstrādāta attiecīgi kodā ierakstītiem makro, ja tika atrastas sakrišanas. Kad sistēma tiks integrēta ar reālu kompilatoru, tā strādās paralēli ar parsētāju un tālāk sistēmas izejas tokenu virkne tiks apstrādāta ar standartiem valodas likumiem.

4.3. Makro konfliktu risināšana

Makro šablonu konflikti var rasties tad, kad daži makro var tikt akceptēti vienlaikus. Tas var notikt gadījumos, ja divas regulārs izteiksmes akceptē līdzīgas virknes. Zemāk tiks aprakstīts, kā tika izvēlēts risināt dažādas konfliktu situācijas.

Reālajā situācijā var rasties 3 konfliktu veidi. Pirmais var rasties gadījumā, kad divas izteiksmes atpazīst vienu un to pašu virkni vienā tvērumā. Otrais var rasties gadījumā, kad jaunajā tvērumā parādās šablons, kas ir līdzīgs jau eksistējošam šablonam no vispārīgāka tvēruma. Trešais var rasties tad, kad viena no izteiksmēm akceptē kādu virkni, bet cita akceptē garāku virkni.

4.3.1. Divu makro konflikts vienā tvērumā.

Gadījumā, ja viena tvēruma ietvaros eksistē divi šabloni, kas dod sakritību ar vienādu garumu, tad tiek ņemtas vērā prioritātes. Tā izteiksme, kas tika ielasīta agrāk būs ar lielāku prioritāti nekā tā, kas ir ielasīta vēlāk. Tātad ja secīgi tiks ielasītas divas izteiksmes $\{id\} \{(\{ \}) \}$ un $\{id\} \{(\{ \{real\} * \}) \}$, tad ielasot virkni $\{id:foo\} \{(\{ \}) \}$ tiks akceptēta pirmā izteiksme. Gadījumā, ja izteiksmes tiks ielasītas pretējā secībā, pirmā izteiksme nekad netiks atpazīta, jo otrā izteiksme pārklāj visas pirmās izteiksmes korektās ieejas.

4.3.2. Divu makro konflikts dažādos tvērumos.

Tvēruma iekšienē strādā tādi paši likumi par izteiksmju prioritātēm - izteiksme, kas bija agrāk ir ar lielāku prioritāti. Bet makro, kas ir specifiski tvērumam ir ar lielāku prioritāti nekā vispārīgāki makro. Tātad, ja pirmajā tvērumā tiks ieviestas makro ar identifikatoriem 1 un 2, bet otrajā tvērumā tiks ieviestas makro ar identifikatoriem 3 un 4, to prioritāšu rinda izskatīsies sekojoši: 3, 4, 1, 2. Pirmie makro ir ar lielāku prioritāti, nekā tie, kas atnāca vēlāk, bet vēlāka tvēruma makro ir ar lielāku prioritāti neka tie, kas atrodami agrākā tvērumā.

4.3.3. Dažādu virkņu garumu konflikts

Prototips strādā pēc mantkārīga (*greedy*) principa - tas akceptē visgarāko iespējamo šablona sakritību. Tātad, ja eksistē divi šabloni $\{int\} \{ , \}$ un $(\{int\} \{ , \}) *$, tad tokenu virkne $\{int:4\} \{ , \}$ $\{int:6\} \{ , \}$ tiks akceptēta ar otru šablonu, neskatoties uz to, ka augstākās prioritātes šablona sakritība tika konstatēta agrāk.

4.4. Realizācijas pamatojums

Galvenais princips uz kura bāzējas prototipa izstrāde ir samazināt apstrādes laiku. Tāad prototipa risinājums tika izveidots tā, lai jebkurā laika momentā šablonu sakrišanu meklēšanai būtu nepieciešams lineārs laiks un tikai viena tokenu virknes apstaigāšana. Šāda pieeja ir izvēlēta ar iedomu, ka makro pievienošana tiks izpildīta tikai vienreiz, un to daudzums būs samērā neliels, bet sakrišanu meklēšana tiks pildīta katrā produkcijā, un, sliktākajā gadījumā, katram tokenam no ieejas plūsmas.

Ielasītā regulārā izteiksme tiek pārsēta un pārveidota nedeterminēta galīgā automātā. Tad šablona nedeterminēts galīgs automāts tiek determinēts un minimizēts. Tāad katrai regulārai izteiksmei tiek izveidots minimāls determinēts automāts, kurš ir optimizēts gan pēc apstaigāšanas laika, gan pēc aizņemtās vietas.

Tālāk, lai nodrošinātu visu šablonu pārbaudi vienlaikus un minimizēt meklēšanas laiku, ir nepieciešams apvienot izveidotos automātus. To var izdarīt dažos veidos. Vienkāršākais no tiek būtu glabāt visus galīgos automātus sarakstā. Pieņemsim, ka ir n šabloni, kurus vajag pārbaudīt. Tad automātu saraksts reprezentē nedeterminētu galīgu automātu ar n ε -zariem no sākuma stāvokļa, katrs no kuriem ved pie sākuma stāvokļa vienam no jau izveidotiem determinētiem automātiem.

Cits veids, kā to varētu izpildīt, ir apvienot visus izveidotos šablonu automātus vienā determinētā galīgā automātā. Tieši šis veids tika izvēlēts šī darba ietvaros lai pēc iespējas samazinātu laiku sakrišanu meklēšanai. Kaut arī automātu apvienošana šādā veidā ir laikietilpīga, tā samazina laika kārtu sakritību meklēšanai.

4.5. Kāpēc neder jau uzrakstītas bibliotēkas

Eksistē diezgan daudz jau izstrādātu bibliotēku, kas apstrādā regulārās izteiksmes. Piemēram, darbs [1] arī lieto automātu teoriju lai paātrinātu apstrādes laiku. Diemžēl tās nav lietojamas tādēļ, ka visi eksistējošie regulāro strādā tieši ar tekstu un simboliem, nevis ar tokeniem un to virknēm. Eksistējošo bibliotēku pārtulkošana uz tokenu apstrādi būtu darbietilpīgāka nekā jaunā moduļa izveide. Šādas bibliotēkas parasti arī piedāvā daudz vairāk funkcionalitātes, nekā ir nepieciešams, piemēram dažādi kontroles simboli, tādi ka \wedge - meklēt no virknes sākuma. Tā kā sistēma tiks izsaukta patvaļīgas tokenu virknes vietas apstrādei, šāda veida kontrole nav nepieciešama. Citu programmētāju rakstītā koda apārstrāde un vienkāršošana aizņemtu vairāk laika nekā tā izveidošana pilnīgi no jauna.

Pieejamas regulāro izteiksmju bibliotēkas arī nepiedāvā iespēju apvienot visus automātus viena vienīgā. Eksistējošie apvienošanas risinājumi, savukārt, pēc apvienošanas neņem vērā, tieši kāda regulārā izteiksme uz doto brīdi ir akceptējosa. Šī darba ietvaros šīs informācijas saglabāšana ir būtiska, jo no tā ir atkarīgs, kāda no transformācijām ir lietota.

4.6. Lietotie algoritmi

Šī apakšnodaļa apraksta algoritmus, kas tika lietoti prototipa realizācijā. Kā jau bija teikts, meklēšanas laika optimizācijai tika izvēlēta pieeja, kur visi regulāro izteiksmju automāti tiek sapludināti kopā.

Visu automātu pārejas pa zariem notiek nevis pa kādu simbolu, bet gan pa attiecīgu tokenu. Regulāro izteiksmju apstrādes gaitā tokeni $\{id:foo\}$ un $\{id\}$ tiek uzskatīti par dažādiem, kaut arī $\{id:foo\}$ ir apakšgadījums tokenam $\{id\}$. Šis fakts tiek iegaumēts tikai sakrišanu meklēšanas gaitā.

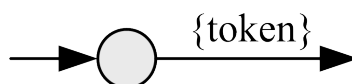
4.6.1. Regulāro izteiksmju pārveidošana nedeterminētā galīgā automātā

FIXME: *Proof that they have the same computational power*

Regulāro izteiksmju translēšana uz nedeterminētu galīgu automātu ir diezgan vienkārša. Lai to paveikt ir nepieciešams pārveidot galvenos regulārās izteiksmes kontroles elementus un automāta gabaliem. Tā kā uz doto brīdi prototips atbalsta tikai ierobežotu regulāro izteiksmju sintaksi, to ir vienkārši izdarīt.

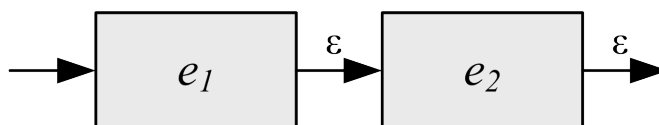
Nedeterminēts galīgs automāts (NGA) veselai regulārai izteiksmei ir izveidots to daļējiem automātiem katrai regulārās izteiksmes daļai. Katram operatoram tiek piekārtota attiecīga konstrukcija. Daļējiem automātiem nav akceptējošu stāvokļu, tiem ir pārejas uz nekurieni, kuras vēlāk tiks lietotas lai savienotu automāta daļas. Pilnīga automāta būvēšanas process beigsies ar akceptējošā stāvokļa pievienošanu palikušajām pārejām. Zemāk tiek parādīti automāti katrai no regulārās izteiksmes iespējamām sastāvdaļām.

Attēlā 5. ir parādīts NGA vienam tokenam *token*.



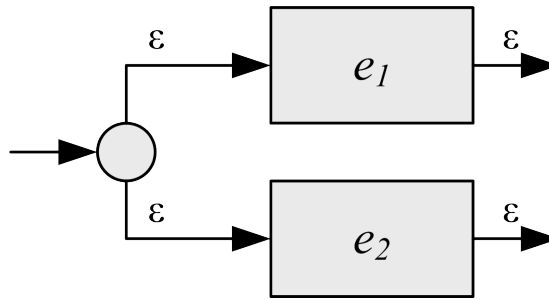
5. att. Automāts vienam tokenam

Attēlā 6. ir parādīts NGA divu automātu konkatenācijai e_1e_2 .



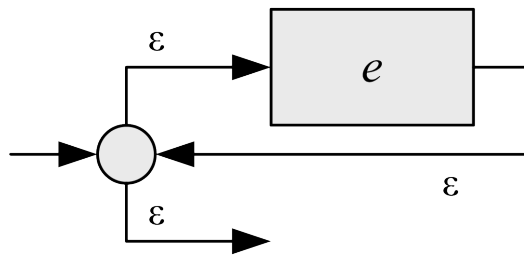
6. att. Automāts divu automātu konkatenācijai

Attēlā 7. ir parādīts NGA izvēlei starp diviem automātiem $e_1|e_2$.



7. att. Automāts izvēlei starp diviem automātiem

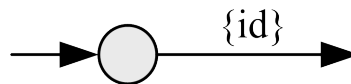
Attēlā 8. ir parādīts NGA priekš konstrukcijas e_1^* .



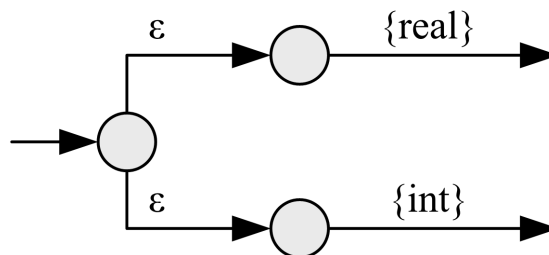
8. att. Automats automātu virknei

Tālāk šie automāti tiek apvienoti vienā, un beigās tiek pievienots akceptējošais stāvoklis.

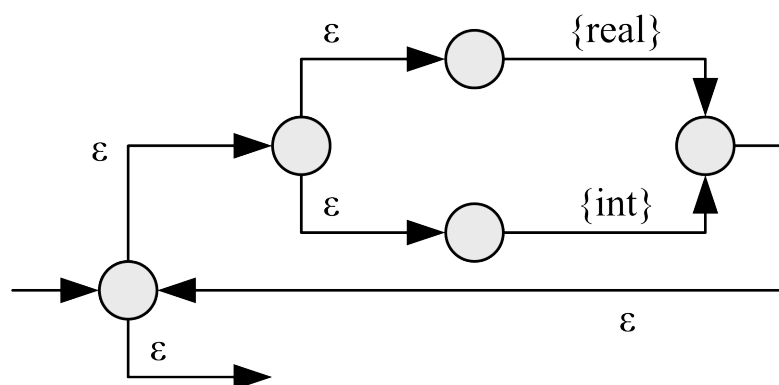
Apskatīsim piemēru - izteiksmi $\{id\} (\{real\} \mid \{int\})^*$. Sākumā tiek izveidoti NGA izteiksmes daļām. Attēls 9. parāda NGA priekš $\{id\}$. Attēls 10. parāda NGA priekš daļas $\{real\} \mid \{int\}$. Attēls 11. parāda NGA priekš $(\{real\} \mid \{int\})^*$.



9. att. Automats tokenam $\{id\}$

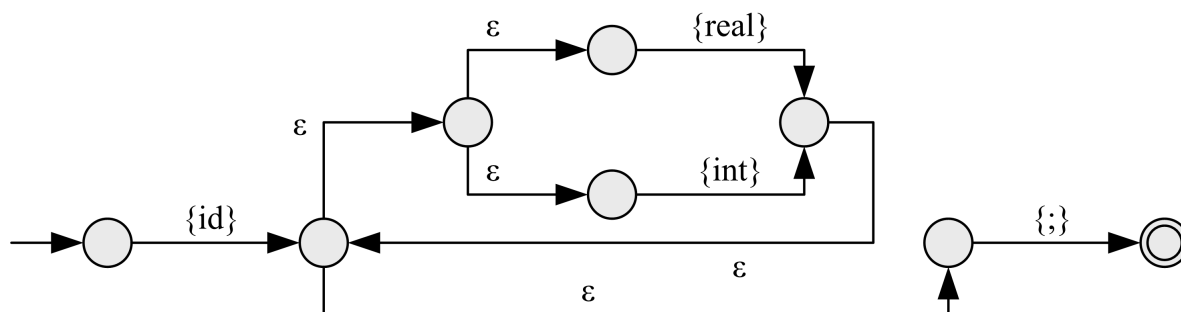


10. att. Automats izteiksmei $\{real\} \mid \{int\}$



11. att. Automats izteiksmei $(\{real\} \mid \{int\})^*$

Tad izveidotie automāti var tikt savienoti un beigās tiem tiek pievienots akceptējošais stāvoklis (attēls 12.).



12. att. Automats izteiksmei $\{id\} (\{real\} \mid \{int\})^*$

Tā tiek izveidots nedeterminēts automāts katrai regulārai izteiksmei. [4], [2]

4.6.2. Determinizācija

Kaut arī daudzām valodām ir vienkāršāk uzbūvēt nedeterminētu galīgu automātu (piemēram, pašām regulārām izteiksmēm tas ir loģiskāk), ir patiesi tas, ka katra valoda var tikt aprakstīta gan ar nedeterminētu, gan ar determinētu galīgu automātu. Turklāt, dzīvī sastopamās situācijās DGA parasti satur tik pat daudz stāvokļu, cik ir NGA. Sliktākajā gadījumā, tomēr var gadīties, ka mazākais iespējamais DGA saturēs m^n stāvokļu (m - ieejas alfabēta elementu skaits), kamēr mazākais NGA saturēs n stāvokļus.¹

Izrādās, ka patiesībā katram NGA eksistē ekvivalents DGA, ko var uzbūvēt ar apakškopu sastādīšanas algoritmu². Vadošā doma šī algoritmā ir tas, ka katrs determinētā galīgā automāta (DGA) stāvoklis ir kādu NGA stāvokļu kopa. Pēc ieejas virknes a_1, a_2, \dots, a_n ielasīšanas DGA atrodas stāvoklī, kas atbilst NGA stāvokļu kopai, kuru var sasniegt apstaigājot virkni a_1, a_2, \dots, a_n .

¹Pieņemsim, ka automāta valoda sastāv no diviem simboliem - 0, 1. Sliktākais gadījums, kad DGA tiešām saturēs 2^n stāvokļus attiecībā pret n NGA stāvokļiem, var rasties tad, kad, piemēram, automāta valodā n -tais simbols no virknes beigām ir 1. Tad DGA būs jāprot atcerēties pēdējos n simbolus. Tā kā ir divi ieejas alfabēta simboli, automātam ir jāatceras visas to dažādas 2^n kombinācijas.

²Pierādījumu tam, ka uzbūvētais DGA tik tiešām akceptē to pašu valodu, ko NGA, sk. [5], teorēma 2.11.

Algoritms 1: NGA transformēšana uz DGA

Ieeja: NGA N .

Izeja: DGA D , kas ir ekvivalents N .

Algoritms: Sākumā algoritms konstruē pāreju tabulu priekš D . Katrs D stāvoklis ir N stāvokļu kopa, tātad tabula tiek konstruēta tā, lai D simulētu vienlaikus visas pārejas, ko var izpildīt N , saņemot kādu ieejas virkni. Lai automāts kļūtu determinēts, ir nepieciešams atbrīvoties no iespējas atrasties dažos stāvokļos vienlaikus. Tātad vajag atbrīvoties no ε -pārejām, un no daudzkārtīgām pārejām no viena stāvokļa pa vienu ieejas simbolu.

Tabulā 1. var redzēt divas funkcijas, kas ir nepieciešamas NGA apstrādes izpildei. Šīs funkcijas no NGA stāvokļiem un pārejām veido jaunas stāvokļu kopas, kuras veidos DGA stāvokļus.

1. tabula. NGA apstaigāšanas funkcijas

Funkcija	Apraksts
$\varepsilon - \text{closure}(T)$	NGA stāvokļu kopa, kas ir sasniedzama lietojot tikai ε -pārejas no visiem stāvokļiem no kopas T .
$\text{move}(T, a)$	NGA stāvokļu kopa, kas ir sasniedzama lietojot pārejas pa simbolu a no visiem stāvokļiem no kopas T .

Ir nepieciešams apstrādāt visas tādas N stāvokļu kopas, kuras ir sasniedzamas, N saņemot kaut kādu ieejas virkni. Indukcijas bāzes pieņēmums ir tas, ka pirms darbības uzsākšanas N var atrasties jebkurā no stāvokļiem, kurus var sasniegt pārejot pa ε bultiņām no N sākuma stāvokļa. Ja s_0 ir N sākuma stāvoklis, D sākuma stāvoklis būs $\varepsilon - \text{closure}(\text{set}(s_0))$. Indukcijai pieņemam, ka N var atrasties T stāvokļu kopā pēc virknes x ielasīšanas. Tad, ja N ielasīs nākamo simbolu a , tad N var pārvietoties jebkura no stāvokļiem $\text{move}(T, a)$. Taču pēc a ielasīšanas var notikt vēl dažas ε -pārejas, tāpēc pēc virknes xa ielasīšanas N var atrasties jebkurā no stāvokļiem $\varepsilon - \text{closure}(\text{move}(T, a))$. Attēls 13. parāda pseido-kodu algoritmam, kā šādā veidā var tikt uzkonstruēti visi DGA stāvokļi un tā pāreju tabula.

Automāta D sākuma stāvoklis ir $\varepsilon - \text{closure}(\text{set}(s_0))$, bet D akceptējošie stāvokļi ir visas tās NGA stāvokļu kopas, kas satur vismaz vienu akceptējošu stāvokli. $Dstates$ ir jauna automāta D stāvokļu saraksts un $Dtran$ ir stāvokļu pāreju tabula.

Sarežģītība: Sarežģītības novērtējums šim algoritmam ir diezgan nepatīkams. Sliktākajā gadījumā tas būs $O(m^n)$, kur n ir NGA stāvokļu daudzums un m ir ieejas alfabēta simbolu skaits. Algoritms var uzģenerēt līdz m^n DGA stāvokļiem, katram no kuram ir m pārejas. Taču parasti tas tā nenotiek un DGA stāvokļu skaits ir līdzīgs NGA stāvokļu skaitam, un algoritma sarežģītība ir $O(n * m)$. [2, 5]

```

initially,  $\varepsilon - \text{closure}(\text{set}(s_0))$  is the only state in  $Dstates$ , and is unmarked
while there is an unmarked state  $S$  in  $Dstates$  do
  mark  $S$ 
  for each available oath  $t$  from  $S$  do
     $U = \varepsilon - \text{closure}(\text{move}(S, a))$ 
    if  $U$  is not in  $Dstates$  then
      add  $U$  as an unmarked state to  $Dstates$ 
    end if
     $Dtran[S, a] = U$ ;
  end for
end while

```

13. att. Automāta determinizēšanas algoritms

4.6.3. Minimizēšana

Izveidotais determinēts galīgs automāts var būt neoptimāls pēc stāvokļu skaita. Bet no šī skaitļa ir atkarīgs tālāko soļu izpildes ātrums. Tāpat ir nepieciešams izveidot automātu, kas atpazīs to pašu valodu un saturēs minimālu iespējamu stāvokļu skaitu.

Var pierādīt, ka katram automātam eksistē ekvivalents minimāls automāts¹. Vēl vairāk, ja eksistē 2 dažādi automāti ar vienādu stāvokļu daudzumu, kas atpazīst vienu un to pašu valodu, tad tie ir vienādi līdz stāvokļu nosaukumiem².

Tālāk teiksim, ka virkne x atšķir stāvokļu s no stāvokļa t tad, kad tikai viens stāvoklis, ko var sasniegt no t un s pa x ir akceptējošs. Tāpat divi stāvokļi ir atšķirami tad, kad eksistē tāda virkne, kas viņus atšķir. Jebkurš akceptējošs stāvoklis ir atšķirams no jebkura neakceptējoša stāvokļa ar tukšu virkni (stāvoklis nevar būt akceptējošs un neakceptējoss vienlaikus). Divi neatšķirami stāvokļi ir ekvivalenti.³

Algoritms 2: DGA minimizēšana

Ieeja: DGA D .

Izeja: Jauns DGA D' , kas ir minimāls un ekvivalents D .

Algoritms: Minimizēšanas algoritma vadošā doma ir sadalīt automātu neatšķiramos stāvokļu grupās. Tas izveido ekvivalentu stāvokļu grupas, kas tālāk var tikt apvienotas vienā, izveidojot minimāla automāta stāvokļus.

Minimizēšanas gaitā automāta stāvokļi tiek sadalīti grupās, ko uz doto brīdi algoritms nevar atšķirt. Jebkuri divi stāvokļi no dažādām grupām ir atšķirami. Katrā nākamajā algoritma iterācijā eksistējošās grupas tiek sadalītas mazākajās grupās, gadījumā, ja kādā grupā parādās atšķirami stāvokļi. Algoritms apstājas līdz ko neviena grupa nevar tikt sadalīta sīkāk.

¹Šī fakta pierādījumu sk. ???

²Tā kā stāvokļu nosaukumi neietekmē automāta darbību, divi automāti tiek saukti par vienādiem līdz pat stāvokļu nosaukumiem, ja viens no tiem var tikt pārveidots otrajā vienkārši pārsaucot to stāvokļus.

³[5] nodaļa 4.4

Pirms algoritms uzsāk darbu, stāvokļi tiek sadalīti divās grupās - akceptējošie stāvokļi un neakceptējošie stāvokļi. Šo grupu stāvokļi ir atšķirami ar tukšu virkni. Tālāk tiek pa vienai apstrādātas grupas no pašreizējā sadalījums. Katrai grupai tiek pārbaudīts, vai tās stāvokļi var tikt atšķirti ar kādu ieejas simbolu - vai kāds no ieejas simboliem noved uz divām vai vairākām dažādam stāvokļu grupām. Ja tādi simboli eksistē, tiek izveidotas jaunas grupas, tādas, ka divi stāvokļi atrodas vienā grupā tad un tikai tad, ja tie aiziet uz vienādām grupām pa vienādiem simboliem. Process ir atkārtots visām pašreizējā sadalījuma grupām, tad atkal jaunam sadalījumam, kamēr neviena no grupām vairs nevar tikt sadalīta.

Attēls 14. parāda minimizēšanas algoritmu pseido-kodā.

```

initially, partitioning  $\Pi$  contains two groups,  $F$  and  $S - F$ , the accepting and nonaccepting
states of  $D$ ,  $\Pi_{new}$  is empty
while  $\Pi_{new}$  is not equal to  $\Pi$  do
     $\Pi_{new} = \Pi$ 
    for each group  $G$  of  $\Pi$  do
        partition  $G$  into subgroups such that two states  $s$  and  $t$  are in the same subgroup if and
        only if for all input symbols  $a$ , states  $s$  and  $t$  have transitions on  $a$  to states in the same group
        of  $\Pi$ 
        replace  $G$  in  $\Pi_{new}$  by the set of all subgroups formed
    end for
end while
 $\Pi_{final} = \Pi$ 

```

14. att. Automāta minimizēšanas algoritms

Tālāk paliek apstrādāt jaunizveidotās stāvokļu grupas izveidojot jaunu determinētu automātu. Lai to izdarītu, no katras sadalījuma Π_{final} grupas tiek izvēlēts grupas pārstāvis. Grupu pārstāvji izveidos jaunus stāvokļus automātam D' . Pārējās komponentes minimālam automātam D' tiks iveidotas sekojoši:

1. Automāta D' sākuma stāvoklis ir pārstāvis tai grupai, kura satur automāta D sākuma stāvokli.
2. Automāta D' akceptējošie stāvokļi ir pārstāvji tām grupām, kuras satur automāta D akceptējošos stāvokļus. Katra no grupām satur vai nu tikai akceptējošus, vai nu tikai neakceptējošus stāvokļus, jo algoritma darba gaitā jaunās grupas tika izveidotas tikai sadalot jau eksistējošas grupas, bet sākuma sadalījums atdalīja šīs stāvokļu klases.
3. Pieņemsim, ka s ir kādas Π_{final} grupas G pārstāvis, un automāts D no stāvokļa s pa ieejas simbolu a pāriet uz stāvokli t . Pieņemsim, ka r ir grupas H pārstāvis, H satur t . Tad automātā D' ir pāreja no stāvokļa s uz stāvokli r pa ieejas simbolu a .

Sarežģītība: Šī algoritma sarežģītība ir $O(n^2 * l)$, kur n ir sākotnējā automāta stāvokļu daudzums un l maksimāls pāreju daudzums no kāda no stāvokļiem ¹.

¹ Patiesībā l būtu jābūt ieejas alfabēta elementu skaitam. Bet šī darba ietvaros ir atļauts veidot regulārās izteiksmes ar pārejām pa tokeniem ar specificētu vērtību, piemēram $\{int : 5\}$. Tātad īstenībā ieejas alfabēts ir bezgalīgs.

4.6.4. Apvienošana

Kā jau bija teikts agrāk, lai samazinātu sakrišanu meklēšanas laiku, tika izvēlēts apvienot visus meklēšanas automātus vienā. Tā kā makro atnāk pa vienam dažādās vietās programmas kodā un var sākt uzreiz tikt lietotas, nav iespējams gaidīt kamēr sakrāsies vairāki automāti apvienošanai. Tikko parādās divi automāti tie tūlīt pat tiek apvienoti vienā sistēmā. Tālāk, kad parādās citi makro, to automāti tiek pievienoti jau eksistējošam.

Algoritms pēc savas būtības ir ļoti līdzīgs determinēšanas algoritmam. Vienīgais uzņēmums ir tas, ka nevienā no automātiem neeksistē ε -pārejas. Tātad tas, no kā vajag atbrīvoties, ir pārejas pa vienu un to pašu simbolu uz diviem dažādiem stāvokļiem. Tas tiek darīts apvienojot divus stāvokļus no dažādiem automātiem.

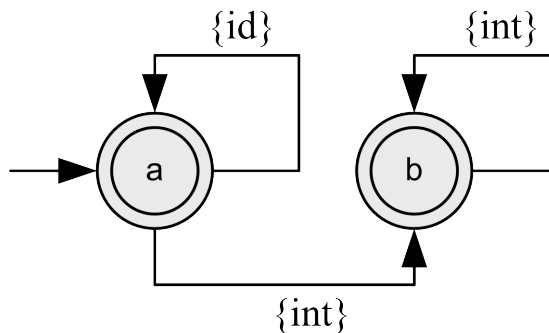
Algoritms 3: Divu DGA apvienošana

Ieeja: DGA D_1 un D_2 .

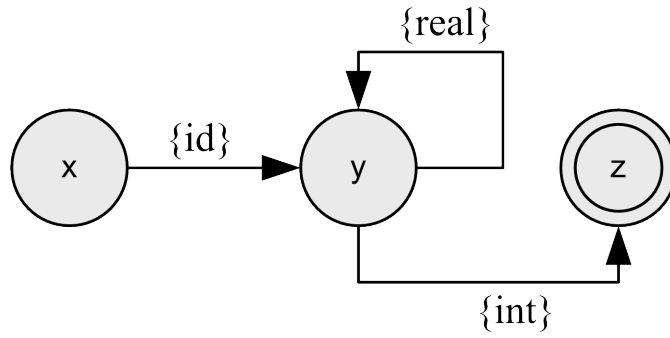
Izeja: Jauns DGA D' , kas apvieno D_1 un D_2 .

Algoritms: Algoritms sāk darbu apvienojot D_1 un D_2 sākuma stāvokļus. Šo stāvokļu kombinācija veido automāta D' sākuma stāvokli. Tālāk algoritms apskata visas iespējamās pārejas no katra no kombinētiem stāvokļiem un apvieno to rezultātus jaunajos stāvokļos.

Apskatīsim piemēru automātu apvienošanai. Attēls 15. parāda determinētu minimālu automātu priekš regulārās izteiksmes $\{id\}^* \{int\}^*$. Tas satur divus stāvokļus, a un b, kuri abi ir akceptējoši. Attēls 16., savukārt, parāda DGA priekš izteiksmes $\{id\} \{real\}^* \{int\}$.

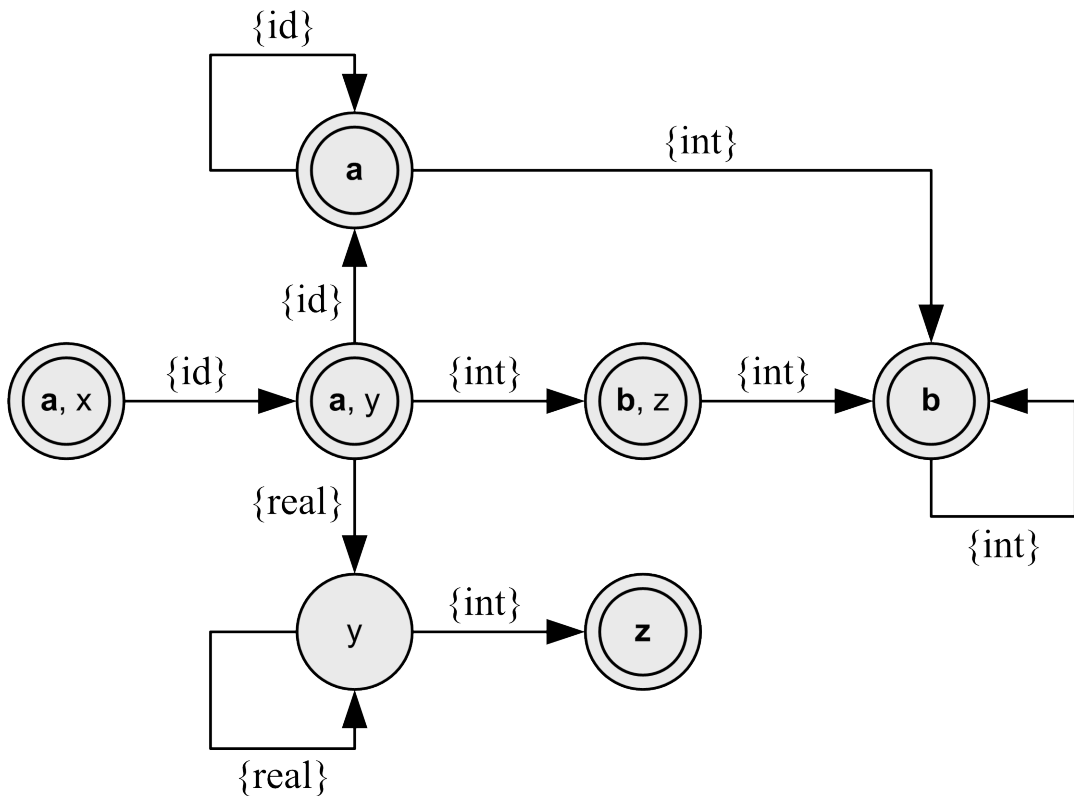


15. att. Automats izteiksmei $\{id\} (\{real\} \mid \{int\})^*$



16. att. Automats izteiksmei $\{id\} (\{real\} \mid \{int\})^*$

To apvienotais automāts ir parādīts attēlā 17..



17. att. Automats izteiksmei $\{id\} (\{real\} \mid \{int\})^*$

Attēls 18. parāda apvienošanas algoritmu pseido-kodā. s_0 ir jaunā automāta D' sākuma stāvoklis, r_0 ir D_1 sākuma stāvoklis un t_0 ir D_2 sākuma stāvoklis. $Dstates$ ir automāta D' stāvokļu saraksts, $Dtran$ ir D' pāreju tabula. Funkcija $move(S, t)$ atgriež tos stāvokļus, uz kuriem var nokļūt no S pa tokenu t . Tā kā parasti S sastāv no diviem stāvokļiem r_i un t_i , tā atgriež pārejas rezultātu no katra no tiem. Rezultāts arī var būt tikai viens stāvoklis, gadījumā, ja no kāda r_i un t_i neeksistē pāreja pa doto tokenu.

Sarežģītība: Šī algoritma sarežģītība ir $O(n * m * l)$, kur n ir pirmā automāta stāvokļu daudzums, m ir otrā automāta stāvokļu daudzums, un l maksimāls pāreju daudzums no katra no stāvokļiem.

```

initially,  $s_0 = (r_0, t_0)$  is the only state in  $Dstates$ , and is unmarked
while there is an unmarked state  $S$  in  $Dstates$  do
    mark  $S$ 
    for each available move  $t$  from  $S$  do  $\triangleright S$  is a combination of some states  $r_i$  and  $t_i$  of  $D_1$ 
    and  $D_2$ , although it might be just a single state from one of the automata.
         $U = move(S, t)$ 
        if  $U$  is not in  $Dstates$  then
            add  $U$  as an unmarked state to  $Dstates$ 
        end if
         $Dtran[S, t] = U$ ;
    end for
end while

```

18. att. Divu automātu apvienošanas algoritms

4.6.5. Sakrišanu meklēšana

Sakrišanu meklēšana NGA slīktākajā gadījumā būs ar sarežģītību $O(n^2 * l)$, kur n ir NGA stāvokļu skaits un l - pārbaudāmās virknes elementu skaits. Tas var notikt, kad visi NGA stāvokļi ir aktīvi vienā laika brīdī, un no katra no tiem eksistē n pārejas pa ieejas simbolu uz visiem automāta stāvokļiem.

Tīrā DGA gadījumā sakrišanu pārbaude katram ieejas elementam ir $O(l)$, kur l ir pārbaudāmās virknes garums.

Diemžēl pilnībā lineāra laika sakrišanu meklēšana nav iespējama tādēļ, ka prototips dod iespēju lietot makro ar tokenu vērtībām. Piemēram, eksistē 2 makro, viens no kuriem gaida tokenu $\{id\}$, un otrs $\{id:foo\}$. Gadījumā, kad sakrišanu meklēšanas procesā parādās tokens $\{id:foo\}$, automātam nav iespējas izsecināt, kurš no ceļiem novedīs pie garākas sakritības. Tādēļ tas iet pa abiem ceļiem vienlaikus, saglabājot abus stāvokļus.

Kaut arī tas ievieš nedeterminētību, tā var parādīties tikai augstāk minētā gadījumā un izveidot ne vairāk ka 2 ceļus vienlaikus. Tātad kaut arī nedeterminētība pastāv, tai ir ļoti maza iespējamība un maza ietekme uz sakrišanu meklēšanas laiku.

4.6.6. Tvērumi

Viena no galvenām šīs sistēmas īpašībām ir iespēja atšķirt programmatūras tvērumus. Ir dažādi veidi, kā var izveidot automātus, kas atšķirtu atsevišķu tvērumu makro. Viens no veidiem varētu būt tāds - katram tvērumam izveidot automātu rindu, kur tvērumam specifiskākie automāti tiks pārbaudīti pirmie. Bet gadījumā, ja ir n iekļautie tvērumi un nepieciešamais makro ir atrodams pirmajā automātā, būs jāizpilda vismaz n meklēšanas, līdz ko pareizais šablons tiks atrasts. Atstājot tikai vienu aktīvu automātu vienā laika brīdī, arī ir dažas iespējas. Varētu visus šablonus likt vienā automātā kopā, un tad pēc izejas no tvēruma attiecīgos šablonus dzēst ārā. Tas nozīmētu, ka katram tvērumam jāatceras makro, kas tika pievienoti, un jāprot dzēst daļu no stāvokļiem ārā no automāta. Bet stāvokļu dzēšana ir laikietilpīga operācija, jo tās izpildīšanai būs nepieciešams apstaigāt visu lielo automātu, dzēšot no tā nevajadzīgos stāvokļus.

Tāpēc tika izvēlēta sekojoša pieeja. Ja prototips darba gaitā sastapās ar tvēruma sākuma

simbolu, tas izveido eksistējošā automāta kopiju un ieliek to kaudzē. Tad automātam tiek pievienotas tvēruma makro. Izejot no tvēruma tā specifiskais automāts tiek izmests ārā un darbs tiek turpināts ar pēdējo automātu no kaudzes, kas atbilst iepriekšējam tvērumam. Šādā veidā jebkurā laika brīdī aktīvs ir tikai viens sakrišanu meklēšanas automāts.

4.7. Izņēmumi

4.7.1. Transformācijas

Izstrādātais prototips nenodrošina ar tokenu virkņu transformācijām, jo tā nav sakrišanu meklēšanas mehānisma uzdevums. Tālākajā sistēmas izstrādes gaitā notiks integrācija ar transformāciju sistēmu vai arī abu mehānismu sapludināšana vienā.

4.7.2. Produkcijas

Prototipā pagaidām nav implementēta apstrādes dalīšana pa gramatikas produkcijām, visas regulārās izteiksmes ir sapludinātas vienā automāta. Regulāro izteiksmju dalīšana pa tipiem tiks ieviesta vēlāk, kad tiks uzsākta integrācija un sadarbība ar reālu parsētāju. Tā varētu tikt implementēta līdzīgi tam, kā tiek realizēti konteksti - pa vienam sapludinātam automātam priekš katra produkcijas tipa.

4.7.3. Tokenu klašu mantošana

Sistēmai nav nekādas informācijas par valodas gramatiku. Tieši tāpēc tokens {real} netiks uztverts ka {expr}, kaut arī racionāls skaitlis ir izteiksme. Tā kā sistēmai jābūt neatkarīgai no valodas gramatikas, šī hierarhija nav iekodējama transformāciju sistēmā. To ir jānodrošina parsētājam, attiecīgi apstrādājot tokenus un apkopojot to nozīmi. Par to arī būs jā rūpējas sistēmas lietotājam rakstot savas makro izteiksmes.

4.7.4. Regulārās izteiksmes daļu grupēšana

Tādēļ, ka aprakstītā pieeja mēģina pēc iespējas minimizēt meklēšanas laiku, tā neatļauj veidot atrasto tokenu grupēšanu, kā tas ir parasti pieņemts regulārās izteiksmēs. Tokenu grupēšana un atpakaļnorādes (*backreferences*) uz tokenu grupām nav atļautas.

Tomēr ja parādīsies nepieciešamība, ir apskatīta arī pieeja, kas dos šādu iespēju. Tas varētu tikt izpildīts, determinējot tikai automāta stāvokļus grupu iekšienē, pēc tam ar ε -pārejām secīgi savienojot grupu automātu akceptējošus un sākuma stāvokļus. Automāts būs determinēts tikai grupu ietvaros, bet tad būs pieejamas atrasto tokenu grupas. Tomēr šāda pieeja neatļaus sapludināt dažus automātus vienā, jo sapludināšanas procedūra sabojās grupēšanu.

4.7.5. Sapludinātā automāta minimizācija

Automātu minimizēšana ir diezgan darbietilpīga operācija, tāpēc tā tiek izpildīta tikai uz atsevišķiem automātiem. Apvienotais visu šablonu automāts var nebūt minimāls, jo dažādu mi-

nimālu automātu apvienošana negarantē šo faktu. Bet tā kā apvienota automāta stāvokļu daudzums var būt ļoti liels, minimizēšanas izpilde var būt neefektīva. Tā kā minimizēšana samazina tikai automāta aizņemto vietu, nevis apstaigāšanas laiku, to šajā gadījumā var izlaist. Sapludinātā automāta minimizēšanu apgrūrina arī tas fakts, ka to stāvokļus vajadzēs atšķirt arī pēc tā fakta, kāds no šabloniem ir akceptēts, nevis tikai pēc tā, vai stāvoklis ir akceptējošs.

4.8. Optimizācijas iespējas

Regulāro izteiksmju optimizēšana uz doto brīdi netiek izpildīta, jo to ir vērts izpildīt uz regulārām izteiksmēm, kas tiks lietoti daudzas reizes. Darba apskatītā situācijā regulārās izteiksmes tiks lietotas tikai vienas programmas ietvaros un to optimizēšanai nav īpašas jēgas. Tomēr bez reāliem piemēriem nevar izšķirt, vai tas izveidos būtisku paātrinājumu šādā konkrētā gadījumā vai nē. Dažas regulāro izteiksmju pārrakstīšanas pieejas ir aprakstītas [6].

Dažreiz divu automātu sapludināšana var izraisīt pārāk lielu stāvokļu daudzumu rašanos. Ja reālajā situācijā tas ietekmēs apstrādes laiku, vai parādīsies nepieciešamība ierobežot automāta aizņemto laiku, var apskatīt iespēju glabāt dažus automātus viena vietā. [6]

Minimizēšanas algoritmu var aizvietot ar citu, ātrāku algoritmu, piemēram, Hopkrofta minimizēšanas algoritmu, kura izpildes laiks ir $O(n \log n l)$, kur n ir automāta stāvokļu daudzums un l ir ieejas alfabēta elementu daudzums.[3]

5. Līdzīgu darbu apskats

FIXME: *related work*

6. Rezultāti

Šī darba ietvaros tika izstrādāts prototips sakrišanu meklēšanas sistēmai. Šīs nodaļas apakšnodaļa 6.1. apskata testēšanas stratēģijas, kas tika lietotas lai nodrošinātu to, ka prototips strādā pareizi. Apakšnodaļa 6.2. apraksta integrēšanas iespējas ar valodas Eq kompilatoru.

6.1. Prototipa testēšana

Prototips tika testēts visā izstrādes laikā. Zemāk tiks aprakstīti automātiski palaižamie testi.

6.1.1. Stresa testēšana

Prototips izstrādes laikā tika testēts ar lieliem automātiski ģenerēto datu apjomiem. Tika ģenerētas patvaļīgas (bet korektas) regulāras izteiksmes ar iekavu un * un | simbolu palīdzību. Katrai regulārai izteiksmei tika izveidota arī simbolu virkne, ko šai izteiksmei jāprot atpazīt. Tad uz vienas un tās pašas regulārās izteiksmes un simbolu virknes tika palaists gan prototips, kas tolaik apstrādāja simbolus, gan Python iebūvētais regulāro izteiksmju apstrādes mehānisms. Vienā piegājienā tika ģenerēti 500 šādi testi. Prototipa beigu izstrādes posmā visi šādi testi tika veiksmīgi izpildīti.

Diemžēl pagaidām šī pieeja netiek implementēta ar tokenu regulārām izteiksmēm, jo nav iespējams pārbaudīt sistēmas darba ekvivalenci ar kādu citu sistēmu. Tāpēc tika veikta intensīva sistēmas testēšana, lai pārbaudītu pēc iespējas vairāk reālajā darbā iespējamo situāciju. Tomēr šādas stresa pārbaudes parādīja ka pats regulāro izteiksmju apstrādes mehānisms strādā korekti.

6.1.2. Sistēmas testēšana

Prototipam tika izveidoti apmēram 20 testi, kas pārbauda to darbību iespējamās situācijās. Tabula 2. parāda konceptuālu testu sadalījumu pa grupām. Dažas grupas pārklājas, jo, piemēram, tvērumu pārbaudošie testi pārbauda arī korektas regulāro izteiksmju prioritātes.

2. tabula. Prototipa testu sadalījums pa grupām

Testa nosaukums	Testa apraksts	Kas tiek pārbaudīts
Testi prioritāšu pārbaudei		
Testi ar dažiem šabloniem vienā tvērumā	Testu gaitā tiek izveidotas dažādas regulāras izteiksmes, kuru aprakstītās valodas pārklājas. Tās tiek ievietotas vienā tvērumā pēc kārtas. Tad tiek pārbaudīts, ka tokenu saraksts tiek akceptēts ar pareizu šablonu atiecībā pret to prioritātēm.	Vai tiek korekti apstrādātas šablonu prioritātes viena tvēruma ietvaros.

2. tabula. Prototipa testu sadalījums pa grupām

Testa nosaukums	Testa apraksts	Kas tiek pārbaudīts
Testi ar dažiem šabloniem vienā tvērumā, kur kāds no šabloniem akceptē garāku virkni	Testu gaitā tiek izveidotas dažādas regulāras izteiksmes, kuru aprakstītās valodas pārklājas. Tās tiek ievietotas vienā tvērumā pēc kārtas. Tad tiek pārbaudīts, ka tiek akceptēta garākā iespējamā tokenu virkne.	Vai tiek korekti apstrādātas šablonu prioritātes viena tvēruma ietvaros.
Testi tokenu vērtībām		
Testi ar tokenu vērtībām	Testu gaitā tiek izveidotas dažādas regulāras izteiksmes ar tokenu vērtībām un bez tām. Tiem tiek padotas dažādas tokenu virknes.	Vai tiek korekti apstrādātas šablonu prioritātes un vērtību sakrišanas.
Testi ar vairākiem pieejamiem stāvokļiem vienlaikus	Testu gaitā tiek izveidotas dažādas regulāras izteiksmes ar tokenu vērtībām. Tiem tiek padotas tokenu virknes ar šādām pašām vērtībām, lai izveidotu situācijas, kad ir pieejami daži stāvokļi vienlaikus.	Vai tiek korekti apstrādātas situācijas, kad parādās nedeterminētība.
Testi tvērumu pārbaudei		
Testi ar tvērumu iekļaušanas dziļumu 1	Testu gaitā tiek izveidotas dažas regulāras izteiksmes, kuru aprakstītās valodas pārklājas. Tās tiek ievietotas nultajā un pirmajā tvērumā. Tad tiek pārbaudīti tokenu saraksti.	Vai tiek korekti apstrādāta tvēruma parādīšanās. Vai tiek korekti apstrādātas šablonu prioritātes starp tvērumiem.
Testi ar tvērumu iekļaušanas dziļumu >1	Testu gaitā tiek izveidotas dažas regulāras izteiksmes, kuru aprakstītās valodas pārklājas. Tās tiek ievietotas dažādos tvērumos ar dziļumu kas ir lielāks par vienu. Tad tiek pārbaudīti tokenu saraksti.	Vai tiek korekti apstrādāti dažādi tvērumu dziļumi. Vai tiek korekti apstrādātas prioritātes starp tvērumiem.
Testi ar izeju no tvēruma	Testu gaitā tiek izveidotas dažas regulāras izteiksmes, kas tiek ieliktas nultajā un pirmajā tvērumā. Tiek pārbaudīts, ka pirmā tvēruma šablons atpazīst tokenu virkni. Tālāk pirmais tvērums tiek pamests un tiek pārbaudīts, ka tā šablons vairs nav aktīvs.	Vai pēc izejas no tvēruma attiecīgie šabloni ir noteikti izdzēsti.

2. tabula. Prototipa testu sadalījums pa grupām

Testa nosaukums	Testa apraksts	Kas tiek pārbaudīts
Testi ar izeju no tvēruma un nākamā tvēruma izveidi	Testu gaitā tiek izveidots 1. līmeņa tvērumš ar regulārām izteiksmēm. Tiek pārbaudīts, ka pirmā tvēruma šablons atpazīst tokenu virknes. Tad šis tvērumš tiek pamests un tiek izveidots jauns pirmā līmeņa tvērumš. Tiek pārbaudīts, ka vecā tvēruma šablons ir izmesti, un ka jaunā tvēruma šablons tiek atpazīti.	Vai pēc izejas no tvēruma attiecīgie šablons ir izdzēsti un pēc ieejas jaunajā tvērumā tiek akceptēti pareizi šablons.
Testi bez šablonu sakritībām		
Testi bez neviena šablona	Testu gaitā tiek izveidota sistēma bez neviena šablona. Tiek pārbaudīts, ka neviena sakritība netiek atrasta.	Vai sistēma korekti apstrādā situāciju, kad nav neviena šablona.
Testi ar šabloniem un datiem kas nesakrīt	Testu gaitā tiek izveidota sistēma ar dažiem šabloniem. Tad tiek padotas tokenu virknes kuras neder nevienam no eksistējošiem šabloniem.	Vai sistēma korekti apstrādā situāciju, kad neviena sakrīšana nav atrasta.

6.2. Prototipa integrēšana Eq

Prototips pagaidām netiek integrēts Eq valodas kompilatorā, bet tas tiek plānots tuvākajā nākotnē. Tā kā prototips tika izstrādāts bāzējoties uz Eq parsētāja īpašībām, to būs viegli integrēt eksistējošā kodā. Tā darbs ir gandrīz neatkarīgs no parsētāja darba un neietekmēs jau eksistējošo programmu darbību.

Prototips piedāvā saskarni lai uzsākt jauna tvēruma apstrādi (funkcija `enter_context()`), lai pamestu tvērumu (funkcija `leave_context()`), lai pievienotu makro (`add_match(regex)`) un lai apstaigātu tokenu virkni meklējot sakrīšanas (`match_stream(stream)`). Integrēšanai prototipu būs jāpapildina ar iespēju padot produkcijas tipu tokenu virkņu apstrādes funkcijām. Prototipa tokenu saņemšanas funkcijas būs jāpārslēdz uz parsētāja piedāvāto saskarni tokenu dabūšanai.

Prototipam ir nepieciešama vienkārša saskarne no eksistējošā parsētāja. Parsētājam jādod pieeja pie tokenu virknes lasīšanas, kā arī jāprot aizvietot atrastās tokenu virknes ar citām virknēm, iespējams, ar citu garumu. Parsētājam arī jāprot pārstartēt tokenu virknes lasīšanu no aizvietotas virknes sākuma. Uz doto brīdi visas šīs iespējas ir implementētas Eq kompilatorā.

Apvienojot parsētāja un sakritību meklēšanas prototipu būs nepieciešams ievietot prototipa funkciju izsaukumus katras produkcijas apstrādes sākumā. Gadījumos, kas parsētājs sastapās ar tokenu, kas identificē makro sākšanos, būs nepieciešams izsaukt regulārās izteiksmes parsēšanas funkciju. Savukārt, kad tiek apstrādātas citas produkcijas, būs nepieciešams izsaukt sakrīšanu meklēšanu. Abu funkciju izsaukumos būs nepieciešams padot arī produkcijas tipu, lai

prototips varētu atšķirt, kādu no automātiem papildināt vai lietot sakrišanu atrašanai. Prototipa funkcijas būs jāizsauc arī tvērumu pārslēgšanu brīdī, lai tas varētu implementēt tvērumu makro prioritāšu sadalīšanu.

7. Secinājumi

Tālāk darbs tik turpināts (šeit var pārfrāzēt Conclusions no raksta melnraksta).

Literatūra

- [1] Re2 regular expression parser. <http://code.google.com/p/re2/>.
- [2] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison Wesley, August 2006.
- [3] Jean Berstel, Luc Boasson, Olivier Carton, and Isabelle Fagnot. Minimization of automata. *CoRR*, abs/1010.5318, 2010.
- [4] Russ Cox. Regular expression matching can be simple and fast. 1 2007.
- [5] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2006.
- [6] Fang Yu, Zhifeng Chen, Yanlei Diao, T. V. Lakshman, and Randy H. Katz. Fast and memory-efficient regular expression matching for deep packet inspection. In *Proceedings of the 2006 ACM/IEEE symposium on Architecture for networking and communications systems*, ANCS '06, pages 93--102, New York, NY, USA, 2006. ACM.