

OTAS - տեքստերի համեմատման խնդիրների լուծման ծրագրային գործիք

Արմեն Վ. Քոչարյան

ՀՀ ԳԱԱ Ինֆորմատիկայի և ավտոմատացման պրոբլեմների ինստիտուտ
e-mail: armkoch@ipia.sci.am

Ամփոփում

Ստեղծված է OTAS (Online Tessellation Automata Syntesator) ծրագրային գործիքը, որը տեքստերի համեմատման խնդիրների ՕՆԱ ճամաշնի նմանակի համար սիմթեզում է տեքստեր համեմատող ծրագրային մոդուլ:

1 Ներածություն

Հաշվողական տեխնիկայի վերջին զարգացմաներն առաջ բերեցին հետաքրքրություն դնած զուգահեռ հաշվարկների մեթոդները և բազմաթիվ բնագավառներում նրանց կիրառություններին: Բջջային ավտոմատները պատկանում են ապահովության ամպլիուդա հաշվողական մոդելների շարքին, որոնք համեստանում են կոմպյուտ հաշվարկների հիմք: Դրանց կառուցվածքը կարելի է դիտարկել որպես զուգահեռ հաշվողական գործիք: Բջջային ավտոմատներն օգտագործվում են բազմաթիվ բնագավառներում տարրեր պրոցեսումների մոդելավորման ժամանակ: Հաշվի առնելով տրամաբանության և ՕՆԱ ավտոմատների էկվիվալենտությունը՝ վերջիններս կարող են բնույթի մեկտեղանի երկրորդ կարգի տրամաբանական բանաձևների [3,4]:

Վերջին տարիներում մեծ ջանքեր են գործադրվում բջջային ավտոմատների ավտոմատ սիմթեզման մեթոդների զարգացման համար [8]: Այդ մեթոդները կառուցման են բջջային ավտոմատներ, որոնք կատարում են կարգընթաց ֆունկցիաների համակարգով նկարագրվող հաշվարկներ:

Մեր մոտեցումը տարրերին է դրամից: Մենք դիտարկում ենք խնդիրների դաս, որի համար ապացուցման նմբ բջջային ավտոմատի գոյուրյունը, և որը կարող է նկարագրվել հատուկ տիպի տրամաբանական բանաձևներով: Այդ դասի համար առաջարկվում է բջջային ավտոմատների սիմթեզման ալգորիթմ:

Այս աշխատանքի նպատակն է լուսաբանել OTAS (Online Tessellation Automata Syntesator) ծրագրային գործիքի հմարափորմայները և օգտվելու կանոնները Տեքստերի Համեմատման խնդիրների(ՏՀԱ) դասի լուծման համար:

ՕՆԱ-ն բջջային ավտոմատ է, որը ներուժվել է K.Inoue-ի և A.Nakamura-ի կողմից և նախատեսված է մատրիցային խզմների ճանաչման համար[1,2]: Այս ավտոմատը միանման վերջավոր ավտոմատների զանգված է, որտեղ անցումների այլիք մեկ անգամ անկյունագծով անցնում է զանգվածի(մատրիցայի) վրայով: Հաշվի առնելով տրամաբանության և ՕՆԱ-ների էկվիվալենտությունը՝ վերջիններս կարող են բնույթի մեկտեղանի երկրորդ կարգի բանաձևների [3]:

Աշխատանքը նվիրված է ՏՀՆ լուծող ՕԽԱ ավտոմատների սիմեբզմանը: Տերստերի համեմատման խնդիրները համեխամում են հայտնի Բառերի Համեմատման խնդիրների (ԲՀՀ) ընդհանրացում: OTAS ծրագրային գործիքով սիմեբզման ՕԽԱ-ն համեխամում է միաշափ բջջային ավտոմատ, որը յուրաքանչյուր քայլում որպես մուտք կարդում է U և V տերստերի միայն երկու հաջորդ տառերը: Վերջում ներկայացված են փորձարկումների արդյունքները ArmCluster անձնական պատագործման համակարգիների կաստերի վրա:

Այժմ բացատրենք տերստերի վերաբերյալ ներմուծված հասկացողորդումները:

Սովորաբար տերստերը հասկացվում են որպես բառերի՝ վերջացվոր A այրութեան շղթաների հաջորդականություն: Գոյություն ունեն բառերի համեմատման բազմաթիվ խնդիրներ(ալգորիթմներ) հիմնված տառերի համեմատման վրա:

Նոյն ճևով իմշաբն բառերի դեպքում(այսինքն տառերի շղթաների), մենք դիտարկում ենք տերստերը որպես բառերի հաջորդականություն: Մենք ավելացրել ենք մեկ նոր սիմվոլ, որը կոչվում է «քամանիչ» և չի պատկանում A այրութեանին՝ «քամանիչ» $\notin A$, որը ծառայում է որպես անջատիչ սիմվոլ և ներառված է $A \cup \{\text{քամանիչ}\}$ այրութեան շղթաների կազմում: $A \cup \{\text{քամանիչ}\} \cup \{\text{այրութեան}\}$ շղթաները վերածվում են A այրութեան բառերի հաջորդականությունների, որոնք բաժանված են «քամանիչ» լրացուցիչ սիմվոլով:

Դիտարկված է տերստերի համեմատման խնդիրների (ՏՀՀ) ՕԽԱ ճանաչելի ենթայաց: Այս դասի յուրաքանչյուր խնդիր տրվում է տերստերի վրա որոշված որոշակի քինար P_{problem} պրեդիկատով, և կայանում է $P_{\text{problem}}(U, V)$ պրեդիկատի գմանատման մեջ կամայական տրված $U = u_1, \dots, u_m$ և $V = v_1, \dots, v_n$ տերստերի համար, որտեղ u_i, v_j բառեր են A այրութեան:

Դրված են որոշակի սահմանափակումներ P_{problem} պրեդիկատի վրա՝ դրանով սահմանափակվում դիտարկվող խնդիրների դասը: Ենթադրվում է, որ գոյություն ունեն բառերի վրա որոշված P_1, \dots, P_r քինար պրեդիկատներ (լոկալ պրեդիկատներ) և $\{0, 1\}^r$ այրութեան մատրիցաների վրա որոշված P ունար պրեդիկատ (գլոբալ պրեդիկատ): $P_{\text{problem}}(U, V)$ պրեդիկատի գմանատումը կարող է կատարվել երկու քայլու, որոնք սահմանվում են ներքում: Այս երկու քայլերի արդյունքում ստացվում է $P_{\text{problem}}(U, V)$ պրեդիկատի արդյունքը:

Առաջին քայլը կայանում է U և V տերստերի համեմատման մեջ: Յուրաքանչյուր u_i բառ համեմատվում է յուրաքանչյուր v_j բառի հետ, այսինքն u_i, v_j յուրաքանչյուր գոյզի համար P_1, \dots, P_r քինար պրեդիկատները պետք է հաշվարկվեն: Այս լոկալ համեմատումների արդյունքում ստացվում է U և V տերստերի $IM_{U, V}$ -ինֆորմացիոն մատրիցան, որտեղ

$$IM_{U, V}(i, j) = (P_1(u_i, v_j), \dots, P_r(u_i, v_j)) \quad i = 1, \dots, m, \quad j = 1, \dots, n:$$

Երկրորդ քայլը կայանում է $P(IM_{U, V})$ գլոբալ պրեդիկատի արժեքը հաշվելու մեջ, որը համեխամում է $P_{\text{problem}}(U, V) = P(IM_{U, V})$ պրեդիկատի արժեքը:

Այս խնդիրը կնշանակենք ՏՀՀ(P_1, \dots, P_r, P):

Հեշտ է նկատել, որ շատ խնդիրներ կարող են ճևակերպվել որպես ՏՀՀ: Շատ ԲՀՀ-ներ, որոնք հայտնի են որպես նմուշի փնտրման, պարզ և զուգահեռ նմուշի փնտրման, պատուհանում նմուշի փնտրման և այլ խնդիրներ, կարող են ընդհանրացվել որպես ՏՀՀ-ներ:

Այժմ բերենք երկու օրինակ:

Օրինակ 1. Նմուշի փնտրման խնդիրը որպես բառերի փնտրման խնդիր կայանում է նրանում, թե արդյո՞ք u բառը հանդիսանում է v բառի ֆակտոր: Նմուշի փնտրման խնդիրը որպես տերստերի փնտրման խնդիր կայանում է նրանում, թե արդյո՞ք U տերստ

համեմատում է V տերսուի ֆակտոր, որտեղ U -ի և V -ի մեջ մտնող բառերը համեմատվում են արբարդությամբ, որպես միավորներ:

Այսուհետեւ որպես լոկալ պրեդիկատ մեմբ կարող ենք օգտագործել P_n մուշտրյան պրեդիկատը: $P_n(M)$ գորպալ պրեդիկատն ընդունում է « λ » արժեք այն և միայն այն դեպքում, եթե M մատրիցային գոյություն ունեն այնպիսի հաջորդական սյուներ, որոնք կազմում են միավոր մատրիցա:

Եթե $V = \text{non-empty word}$ և $U = \text{empty word}$, որտեղ $U = \text{empty word}$, որտեղ պրոբել համեմատում է բաժանի, ապա տերսուի համար նմուշի փնտրման խնդիրը ունի բացասական լուծում:

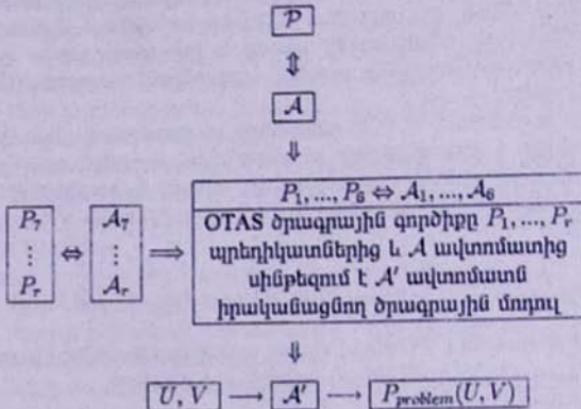
Եթե մեմբ դիտարկենք V -ն և U -ն որպես բառեր այնպիսի այրութեանում, որը պարունակում է բացակայություն, ապա ապա տառ, ապա բառերի համար նմուշի փնտրման խնդիրը ունի դրական լուծում:

Օրինակ 2. Նմուշի գուգահեռ փնտրում:

Հիցոր տրված են $U = u_1, \dots, u_n$ և $V = v_1, \dots, v_m$ երկու տերսուերը: Խնդիրը կայանում է նրանում, թե գոյություն ունի՞ արյուր V -ում այնպիսի v_{j_1}, \dots, v_{j_m} բառերի ենթահօրդակամություն, որ $u_i = v_{j_i}$ բոլոր $i = 1, \dots, m$ համար: Այսուհետեւ մտնութեան կարող ենք օգտագործել որպես լոկալ պրեդիկատ P_n մուշտրյան պրեդիկատը: $P_n(M)$ գորպալ պրեդիկատն ընդունում է 1 արժեք այն և միայն այն դեպքում, եթե M մատրիցայի յուրաքանչյուր տող պարունակում է նվազագույնը մեկ հար 1:

Ստեղծված է OTAS (Online Tessellation Automata Syntesator) ծրագրային գործիք, որը տերսուի համեմատման խնդիրների ՕնU [7] ճանաչելի ենթադասի համար սիմբեզում է տերսուի համեմատող ծրագրային մոդուլ:

Նկար 1-ում ցոյց է տրված OTAS ծրագրային գործիքի կառուցվածքային պատկերը:



Նկար 1. OTAS ծրագրային գործիքի աշխատանքի կառուցվածքային պատկեր:

OTAS ծրագրային գործիքը բավարարում է մախապես դրվագ հետևյալ պահանջներին.

1. այն շատ պարզ է օգտագործման համար,

2. OTAS-ն ունի ներդրված բառերի համեմատման լոկալ պրեդիկատներ,

3. ապահովում է բառերի համեմատման լոկալ պրեդիկատներ ավելացնելու հմարակորդություն,

4. որպես գլորալ պյուղիկատ այն ընդունում է իմֆորմացիոն մատրիցայի վրա աշխատող խճանկարային ավտոմատի անցումների ֆունկցիան,

5. որպես աշխատանքի արդյունք վերադառնում է ավտոմատ կերպով սիմեթրիկ սիմեթրիկ ծրագրային մոդուլ, որը նոյնականացնում է մուտքագրված խճանկարային ավտոմատի աշխատանքը համեմատվող տեքստերի վրա,

6. ավտոմատ կերպով սիմեթրիկ սիմեթրիկ ծրագրային մոդուլ աշխատում է զծային ժամանակում,

7. իրականացված են OTAS ծրագրային գործիքի հաջորդական և գուգահեռ տարրերակմբ: Հաջորդական տարրերակում գեներացվում է այնպիսի ծրագրային մոդուլ, որն աշխատում է մեկ պրոցեսորի վրա որպես հաջորդական ալգորիթմ: Չուզահեռ տարրերակում սիմեթրիկ սիմեթրիկ պրոցեսոր մոդուլ, որն աշխատում է բազմապրոցեսորային կլաստերի վրա որպես գուգահեռ ալգորիթմ:

2 Սահմանումներ

Ենթադրենք A -ն վերջավոր ոչ դատարկ այրութեն է: A այրութենում $w = w_1w_2...w_k$, $w_i \in A$, $i = 1...k$ տառերի հաջորդականությունը կանվանենք $ρառ$: $w \in A^+$, որտեղ $A^+ -$ ով կնշանակենք A այրութենումը բոլոր բառերի բազմությունը՝ բացի դատարկ բառից:

$A^+ \cup \{*\}$ այրութենում $t_1, ..., t_n$ բառերից կազմված $T = *t_1 * ... * t_n *$ տիպի հաջորդականությունը կանվանենք $տեքստ$, որտեղ $t_i \in A^+, i = 1, ..., n$ և $* \notin A$:

T տեքստի երկարությունը կանվանենք $|T| = n + 1 + \sum |t_i|$, որտեղ $|t_i|$ -ն t_i բառի երկարությունն է:

Ենթադրենք U -ն և V -ն տեքստեր են A այրութենում:

$$U = *u_1 * ... * u_m *, \text{ որտեղ } u_1 = u_{11}u_{12}...u_{1k_1}, \dots, u_m = u_{m1}u_{m2}...u_{mk_m} \quad u_{ij} \in A:$$

$$V = *v_1 * ... * v_n *, \text{ որտեղ } v_1 = v_{11}v_{12}...v_{1l_1}, \dots, v_n = v_{n1}v_{n2}...v_{nl_n} \quad v_{ij} \in A:$$

Ենթադրենք ֆիբուլան են $P_1, ..., P_r$ բինար պրեիկատները՝ որոշված A այրութենում բառերի գույքերի վրա: Կամայական (U, V) տեքստերի գույզին կամապատասխանեցնենք (m, n) մատրիցա և կնշանակենք $IM_{U,V}$, որի էլեմենտները r բիտանց վեկտորներ են: $IM_{U,V}(i, j) = (P_1(u_i, v_j), ..., P_r(u_i, v_j))$, $1 \leq i \leq m, 1 \leq j \leq n$: $IM_{U,V}$ մատրիցան կանվանենք (U, V) գույզի իմֆորմացիոն մատրիցա:

$(P_1(u_1, v_1), ..., P_r(u_1, v_1))$	$(P_1(u_1, v_n), ..., P_r(u_1, v_n))$
$(P_1(u_2, v_1), ..., P_r(u_2, v_1))$	$(P_1(u_2, v_n), ..., P_r(u_2, v_n))$
\vdots		\vdots
$(P_1(u_m, v_1), ..., P_r(u_m, v_1))$	$(P_1(u_m, v_n), ..., P_r(u_m, v_n))$

Նկար 2. (U, V) տեքստերի իմֆորմացիոն մատրիցա

(U, V) տեքստերի գույզի մատրիցային կող կանվանենք $(կնշանակենք $MC_{U,V}$)$ $A \cup \{*\}$ այրութենում այնպիսի ($|U|, |V|$) մատրիցան, որի առաջին սյունակը գրված է U տեքստը, առաջին տողում՝ V տեքստը, իսկ մնացած էլեմենտները կամայական սիմվոլներ են վերցված A այրութենից:

*	v_{11}	...	v_{1k_1}	*	v_{21}	...	v_{2k_2}	*	...	*	v_{m1}	...	v_{mk_m}	*
u_{11}														
\vdots														
u_{1k_1}														
*														
u_{21}														
\vdots														
u_{2k_2}														
*														
\vdots														
*														
u_{m1}														
\vdots														
u_{mk_m}														
*														

Նկար 3. (U, V) տերստերի մատրիցային կող

Օրինակ 3. Ենթադրենք $U = *I * \text{match} * \text{texts}*$ և $V = *Now * I * \text{match} * \text{two} * \text{texts}*$: Դիցուք P_u -ը բառերի հավասարության պրեդիկատն է:

$$P_u(u, v) = \begin{cases} 1, & \text{եթե } u = v \\ 0, & \text{եթե } u \neq v \end{cases}$$

U և V տերստերի մատրիցային կողը հետևյալն է:

*	N	e	w	*	I	*	m	a	t	c	h	*	t	w	o	*	t	e	x	t	s	*
I	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
*	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
m	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
a	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
t	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
c	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
h	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
*	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
t	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
e	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
x	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
i	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
s	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o
*	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o

Նկար 4. Օրինակ 3-ում նշված տերստերի մատրիցային կող

Բառերի հավասարության պրեդիկատով U և V տերստերի իմֆորմացիոն մատրիցան

հետևյալն է՝

$$IM_{U,V} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Նկար 5. Օրինակ 3-ում նշված տեքստերի ինֆորմացիոն մատրիցան
 P_+ պրեդիկատով

3 Ծրագրի աշխատանքի նկարագրությունը

Դիցուք օգտվողի խնդիրը կայանում է հետևյալում՝ անհրաժեշտ է ծրագիր, որը U, V տեքստերի գույզի համար հաշվի $\mathcal{P}(U, V)$ պրեդիկատի արժեքը

$$\mathcal{P}(U, V) = \begin{cases} 1, & \text{եթե } \text{տեքստերը } \text{համեմատելի են} \\ 0, & \text{հակառակ դեպքում:} \end{cases}$$

Դիցուք այդ պրեդիկատն արտահայտվում է $(S_1, S_2, P_1, \dots, P_r)$ սիգմատորպայում EMSO տրամարանության մեջ, որտեղ P_1, \dots, P_r բառերի համեմատման պրեդիկատները ՕԽԱ ճանաչելի են:

Հայտնի է, որ մատրացին լեզուների վրա EMSO տրամարանության բանաձևի իրազործվածի համարժեք է ՕԽԱ ավտոմատով ճանաչվությանը: Ուստի համեմատման \mathcal{P} պրեդիկատը կարելի է բերել մատրիցաների վրա համարժեք ՕԽԱ ավտոմատի [3]:

Եթե բավարարված է վերը նշված պահանջը, ապա կարելի է խնդիրը լուծել OTAS ծրագրային գործիքի միջոցով:

Օրինակ 4. Դիտարկենք հետևյալ խնդիրը. հանդիսանում է՝ արդյոք U տեքստը V տեքստի հերատեքստը: Օգտագործենով OTAS-ը՝ կարող ենք ստանալ ծրագրային մոդուլ, որը կլուծի խնդիրը տեքստերի կամայական U, V գույզի համար:

Դիտարկենք կոնկրետ օրինակ: Դիցուք $U = *I * \text{match} * \text{texts}* \text{ և } V = *Now * I * \text{match} * two * \text{texts} * \text{using} * \text{OTAS}*$:

Բառերի հավասարության պրեդիկատով U և V տեքստերի ինֆորմացիոն մատրիցան հետևյալն է՝

$$U \quad V \\ \begin{array}{|c|c|c|c|c|c|c|} \hline & 0 & 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

Նկար 6. Օրինակ 4-ում նշված տեքստերի ինֆորմացիոն մատրիցան
 P_+ պրեդիկատով

Որպեսզի U տեքստը լինի V -ի նմբատեքստ, անհրաժեշտ է և բավարար, որպեսզի ինֆորմացիոն մատրիցայի յուրաքանչյուր տողում հնարավոր լինի ընտրել առնվազն մեկ հատ 1 պարունակող վանդակ այնպիսին, որ յուրաքանչյուր տողում ընտրված վանդակը

գունդի նախարդ տողում ընտրվածից այ:

A -ով նշանակենք $\{0, 1\}$ այլորենում մատրիցաների վրա աշխատող այն ՕԽԱ ավտոմատը, որը լուծում է ներառական գունդու խնդիրը:

A ՕԽԱ ավտոմատն ունի հետևյալ տեսքը

$A = (\Sigma, Q, q_0, F, \delta)$, որտեղ

$\Sigma = \{0, 1\} \cup a \in \Sigma$;

$Q = \{q_0, q_1, q_2, q_3, q_4\}$;

q_0 -ն սկզբանական վիճակն է;

$F = q_2, q_3$:

Ավտոմատի անցումների ժ ֆունկցիան ունի հետևյալ տեսքը՝

$\delta(q_0, q_0, 0) = q_1, \quad \delta(q_1, q_0, 0) = q_1, \quad \delta(q_2, q_0, 0) = q_3, \quad \delta(q_3, q_0, 0) = q_3,$

$\delta(q_0, q_0, 1) = q_2, \quad \delta(q_1, q_0, 1) = q_2, \quad \delta(q_2, q_0, 1) = q_3, \quad \delta(q_3, q_0, 1) = q_3,$

$\delta(q_0, q_1, 0) = q_1, \quad \delta(q_1, q_1, 0) = q_1, \quad \delta(q_1, q_2, 0) = q_1, \quad \delta(q_1, q_3, 0) = q_1,$

$\delta(q_2, q_1, 1) = q_1, \quad \delta(q_1, q_1, 1) = q_1, \quad \delta(q_1, q_2, 1) = q_1, \quad \delta(q_1, q_3, 1) = q_2$:

OTAS-ը ստուարով A ավտոմատի անցումների ֆունկցիան իրականացնող ծրագրային կողը որպես մուտք սիմեբում է ներահաջորդականություն գունդու խնդիրի լուծման ծրագրային մոդուլ: Ստացված ծրագրի միջոցով հնարավոր է լուծել խնդիրը U, V տերստերի կամայական գույքի համար:

OTAS ծրագրային գործիքը խնդիրը լուծելու համար անհրաժեշտ է.

1) Անդրյալած համեմատման սլեյնիկատմերից ընտրել P_{eq} .

2 C++ ծրագրալորման լեզվով նկարագրել իմփորմացիոն մատրիցայի վրա աշխատող խճանկարային A ավտոմատի անցումների ֆունկցիան:

3) մուտքազնի համեմատվող տերստերը պարունակող ֆայլերի հասցեները,

Որպես աշխատանքի արյունը OTAS ծրագրային գործիքը վերաբերենում է աշխատող ծրագրային մոդուլ, որը մողելավորում է A ավտոմատի աշխատանքը համեմատվող տերստերի մատրիցային կողի վրա, այսինքն աշխատում է անմիջապես մուտքային տերստերի վրա: Օգտագործելով ստացված ծրագրային մոդուլը, կարելի է լուծել ներառական գունդու խնդիրը տերստերի կամայական գույքի համար: Օ

Հասուն նշենք, որ OTAS ծրագրային գործիքում որած ներդրված լոկալ պրեյլիկատմեր օգտագործվում են $P_{eq}, P_{eq}, P_{fact}, P_{seq}, P_{req}, P_{perm}$ բառերի համեմատման սլեյնիկատմերը, որոնք բոլոր են տախում լուծել բազմաթիվ Բառերի Համեմտման դասական Խնդիրները(ԲՃՆ) առանց որևէ նոր պրեյլիկատ ավելացնելու: Բազմաթիվ ԲՃՆ-ներ կարելի են ներկայացնել որպես $S\Box(P_1, \mathcal{P})$, որտեղ P_1 -ը ներդրված պրեյլիկատմերից որևէ մեկն է, իսկ \mathcal{P} -ը բոլոր մատրիցաները ճանաչող համապատանի պրեյլիկատը: Խնդիրը լուծելու համար խնդրականիք ուժիմում ընտրվում է համապատասխան պրեյլիկատը, իսկ P_1 -ը ներկայացնելու անհրաժեշտություն չկա: Որպես ծրագրի արյունը պատր է պահանջները ընտրված P_1 պրեյլիկատի առժեքը: Համեմատվող և և բառերը տրվում են մուտքային ֆայլերի մեջ, որոնք վերջանանում են * նիշով:

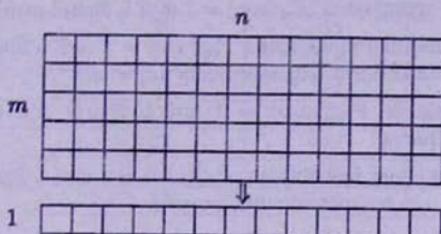
A -ով նշանակենք տերստերի համեմատման խնդիրը լուծող սիմերզված ՕԽԱ ավտոմատը:

Ծրագիրը իրականացված է երկու տարրերակով:

1.Հաջորդական աշխատող մողովի սիմերզման ժամանակ տվյալների երկափ գանգվածի վրա աշխատող խճանկարային A_1 ավտոմատի աշխատանքը մողելավորված է տվյալների միաշափ զանգվածի վրա աշխատող A_2 ավտոմատով: Ընդ որում, եթե A_1 ավտոմատն աշխատում է (m, n) չափերով տվյալների երկափ զանգվածի վրա, ապա A_2 ավտոմատն աշխատում է $(1, n)$ չափերով տվյալների միաշափ զանգվածի վրա և նրա

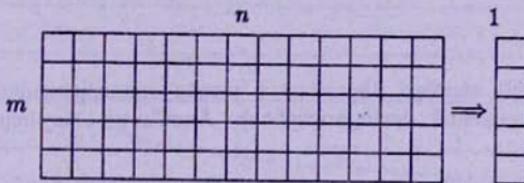
յուրաքանչյուր էլեմենտում պահված է նույն ծավալի իմֆորմացիա, ինչ որ A_1 ավտոմատի աշխատանքի դեպքում: Որպես բջջային ավտոմատ A_2 ավտոմատի աշխատանքի բայերի քանակը $m + n + 1$ է: Էնջան նաև A_2 ավտոմատի աշխատանքի դեպքում միաշափ գանգվածի յուրաքանչյուր բջջում աշխատում է այն նույն էլեմենտուր ավտոմատը, ինչ որ A_1 ավտոմատի աշխատանքի դեպքում: Մեկ բջջում աշխատող էլեմենտուր ավտոմատներն աշխատում են ճիշտ այնքան անգամ, ինչքան նրանք աշխատում են A_1 ավտոմատում:

Նկար 7-ում ցոյց է տրված A_1 ավտոմատից A_2 ավտոմատին անցնելու դեպքում ստացվող հիշողության խնայողությունը, եթե (m, n) չափերով երկշափ մատրիցայի դեպքում օգտագործվում է $(1, n)$ չափերով միաշափ գանգված:



Նկար 7: ОТА ավտոմատի հաջորդական մոդելավորման սխեման

2. Չուզահեռ աշխատող մոդուլ սիմքեզման դեպքում A_1 ավտոմատի աշխատանքը մոդելավորվում է միաշափ սիստեմիկ ալգորիթմի աշխատանքով, որի երկարությունը m է, իսկ քայլերի քանակը՝ $m + n + 1$: Սիստեմիկ ալգորիթմի յուրաքանչյուր բջջում պահվում է նույն քանակի իմֆորմացիա և աշխատում է այն նույն էլեմենտուր ավտոմատը, ինչ որ A_1 ավտոմատի աշխատանքի դեպքում: Սիստեմիկ ալգորիթմի աշխատանքը գուգահեռացվում և մոդելավորվում է կլաստերային միջավայրում:



Նկար 8: ОТА ավտոմատի գուգահեռ մոդելավորման սխեման

Նկար 8-ում սխեմատիկ կերպով ցոյց է տրված A_1 ավտոմատից սիստեմիկ ալգորիթմի անցնելու դեպքում մատրիցայի չափերի փոփոխությունը:

Որպես ավտոմատի մուտքային այրութեն և համեմատվող տեքստերի այրութեն ընտրված է C++ ծրագրավորման լեզվի «char» տիպի սիմվոլների քազմությունը՝ քայլի ֆայլի վերջի սիմվոլից և '*' սիմվոլից, որը համեմատում է անջատիչ:

Մուտքային տվյալների ֆայլերը պետք է պարունակեն համեմատվող քառերի հաջորդականություններ՝ քաժանված '*' սիմվոլով, և ֆայլի վերջին սիմվոլը նույնպես պետք է լինի քաժանիչ սիմվոլ: Չի բույլատրվում հաջորդականության մեջ նշել դատարկ

բառ:

OTAS ծրագրային գործիքն օգտվողին հնարավորություն է տալիս նոր ծրագրային մոդուլի սիրթեզման ժամանակ որպես համեմատման պրեդիկատ ինտերակտիվ միջավայրում ընտրել ծրագրում Ընթրոված P_{eq} , P_{eg} , P_{fact} , P_{seq} , P_{perp} , P_{perm} վեց բառերի համեմատման լրկալ պրեդիկատներից մեկը կամ մի բանից միաժամանակ: Ինչպես նաև օգտվողը հնարավորություն ունի անհրաժեշտության դեպքում ավելացնել կամայական ՕԿԱ ճանաչելի բառերի համեմատման պրեդիկատ, որի համար անհրաժեշտ է մուտքազերել պրեդիկատի անոնց և պրեդիկատին համարժեք խճանկարային պկուսնական անցումային ֆունկցիան՝ նկարագրված C++ ծրագրավորման լեզվով:

Ծրագրում ներդրված լրկալ պրեդիկատներն են՝

1. Հավասարության պրեդիկատ՝ $P_{eq}(u, v) = 1$ այն և միայն այն դեպքում, եթե $u = v$:
2. Հավասար երկարության պրեդիկատ $P_{eq}(u, v) = 1$ այն և միայն այն դեպքում, եթե համեմատվող բառերը ունեն նույն երկարությունը՝ $|u| = |v|$:
3. Ֆակտոր պրեդիկատ՝ $P_{fact}(u, v) = 1$ այն և միայն այն դեպքում, եթե u բառը համեյանում է v -ի ֆակտորը:
4. Հաջորդական ենթարար պրեդիկատ՝ $P_{sep}(u, v) = 1$ այն և միայն այն դեպքում, եթե u բառը համեյանում է v -ի հաջորդական ենթարար:
5. Չուզահեռ ենթարար պրեդիկատ՝ $P_{sep}(u, v) = 1$ այն և միայն այն դեպքում, եթե u բառը համեյանում է v -ի զուզահեռ ենթարար:
6. Տեղադրության պրեդիկատ՝ $P_{perm}(u, v) = 1$ այն և միայն այն դեպքում, եթե u բառը համեյանում է v -ի տեղադրությունը, այսինքն բառերն ունեն նույն երկարությունը և հավասար բանակով այբուբենի յուրաքանչյուր տառից:

4 Փորձարկումներ

OTAS ծրագրային գործիքը ներդրված է բարձր արտադրողականություն ունեցող անձնական օգտագործման կոմպյուտերների ArmCluster հաշվողական կլաստերի միջավայրում:

Գեներացված զուգահեռ ծրագրային մոդուլի գործակության էֆեկտիվությունը ստուգելու համար ArmCluster բազմապրոցեսորային հաշվողական համակարգի վրա կատարվել են հետևյալ խնդրի փորձարկումները: Տեքստում համեյապում է՝ արդյոք նմուշի որևէ բառ երկու անգամ հաջորդարար:

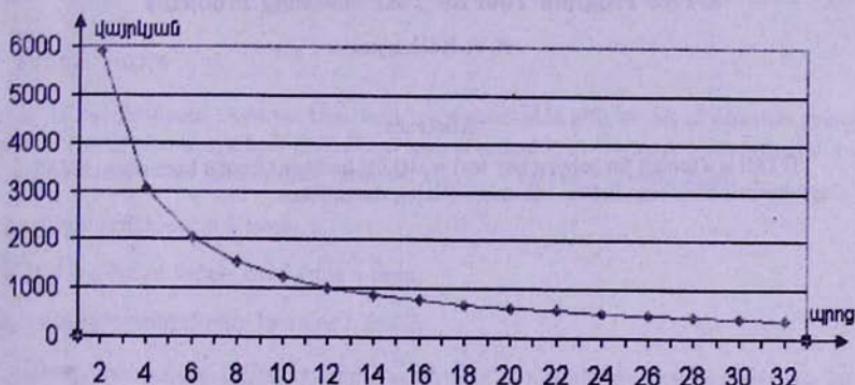
Փորձարկումները կատարվել են տեքստի և նմուշի երեք տարբեր երկարությունների համար 2,4,...,32 պրոցեսորների վրա:

Ալյուսակ 1-ում բերված են OTAS ծրագրային գործիքի միջոցով գեներացված ծրագրային մոդուլի աշխատանքի ժամանակը կախված պրոցեսորների բանակից: Եթեր փորձարկումների դեպքում որպես համեմատված տեքստ և նմուշ վերցվել են նոյն հաշորդականությունները, որոնք առաջի փորձում ունեն 26510 սիմվոլ, երկրորդում 122910 սիմվոլ և երրորդում՝ 1207410 սիմվոլ:

Աղյուսակ 1

Պ.Զ./Տ.Ե	26510	122910	1207410
2	269	5895	568222
4	139	3054	295375
6	94	2061	199073
8	71	1538	150334
10	57	1226	120938
12	48	1020	100795
14	41	876	86642
16	35	768	75760
18	32	682	67542
20	29	615	60864
22	27	559	55331
24	24	512	50759
26	22	473	46866
28	21	440	43572
30	19	410	40606
32	18	385	38069

Նկար 9-ում ցույց է տրված 122910 սիմվոլ պարումակող տեքստերով 2-ից 32 պրոցեսորների վրա կատարված փորձարկումների դեպքում ծախսված ժամանակի կախվածությունը օգտագործվող պրոցեսորների քանակից:



Նկար 9: Զուգահեռ ծրագրի աշխատանքի ժամանակի կախվածությունը պրոցեսորների քանակից

Փորձարկումների արդյունքները հիմնավորում են OTAS ծրագրային գործիքի սիմբեզած ծրագրային մոդուլի համարյա իդեալական զուգահեռացումը: Սիմբեզած ծրագրային մոդուլը աշխատանքի ընդացքում յուրաքանչյուր վայրկյանում մոտավորապես վերամշակում է $1,5 \times 10^4 \times N$ սիմվոլ մոտքային ֆայլներից, որտեղ N -ն աշխատող պրոցեսորների քանակն է:

Օգտագործված զանկանություն

- [1] K.Inoue and A.Nakamura. Some properties of two-dimensional on-line tessellation acceptors. *Information Sciences*, vol.13, p.95-121, 1977.
- [2] K.Inoue and A.Nakamura. A Survey of two-dimensional automata theory. In Proc. 5th Int. Meeting of Young Computer Scientists.J.Dasson and J.Kelemen (Eds.), p. 72-91. Lecture Notes in Computer Science 381, Springer-Verlag, Berlin, 1990.
- [3] D.Giammarresi and A.Restivo. Two-dimensional languages. In *Handbook of Formal Language Theory*, v.3, Springer-Verlag, N.Y.,1996.
- [4] D.Giammarresi, A.Restivo, S.Siebert and W.Thomas. Monadic second order logic over rectangular pictures and recognizability by tiling systems. *Information and computation*. Vol.125, No. 1, p. 32-45, 1996.
- [5] L.Boasson, P. Cegielski, I. Guessarian, Y. Matiyasevich. Window accumulated subsequence matching is linear. *Annals of Pure and Applied Logic* Vol. 113(2001),pp.59-80.
- [6] K. V. Shahbazyan, Yu. H. Shoukourian. The Finite-state Recognizability of sequences of Integers. *Transactions of the Institute for Informatic and Automation Problems of NAS RA*. Vol 27,(2006), 151-173
- [7] A. V. Kocharyan, K. V. Shahbazyan. On-line Tessellation automata as a text Matching Problems Solver. *Transactions of the Institute for Informatic and Automation Problems of NAS RA*. Vol 27,(2006), 54-62.
- [8] C.T. Djamegni and M. Tchuente. A New Dynamic Programming Algorithm on two dimensional arrays. *Parallel Processing Letters*, 10(1)(2000) 15-27.

OTAS Program Tool for Text Matching Problems

A.V. Kocharyan

Abstract

OTAS is a toolkit for solving any text matching problems from a fixed class. OTAS synthesizes online tessellation automata solving the problem.