

Óíèååðñàëüíáy ñèñòåìà ñèíòàéñè÷åñêíé ðàçìåòèè òåéñòà ObjectATE

Àâòîð Àëåéñåé Éäíðååé÷ Çiáíèí

27.06.2008 á.

Íñëäääíáá íáñëäääíéá 25.07.2008 á.

Òåçèñû á ôîðìàòå DOC (167.5 kB 2008-07-15 12:23:14) Òåçèñû á ôîðìàòå PDF (225.08 kB 2008-07-15 12:19:13)

Ñèñòåìà

Ðàçìåòèè òåéñòà ObjectATE (îò Object-oriented ancient text editor) ðàçðàáàòúâååòñý è èñïíëüçóåòñý á ìòäåéå ëëíäåéñòè÷åñé èñòî÷íééâåååíéý ÈÐÙ ÐÀÍ ñ 2006 áíäà [Çiáíèí è äð. 2006]. Òåéóùåy âåðñèý ñèñòåìû ðåàéèçíàáíá íá íæòòîðíå Microsoft .NET Framework 2.0 ñ èñïíëüçíàáíéåí ðåéýöèíííé áàçû äàííûõ Microsoft Access.

Â

Íñííâá ñèñòåìû ëåæèò íáúåéòíí-íðéäíóèðíåáííûé íäðöíá, ñèðîéí íðéäíáíýáíûé á ïðíäðåííéé. Áåñü ðåçìå-åííûé áíéòíáíó íðåñòååëýåòñý êåé íááíð íáúåéòíá. Íðíöåññ ðåçìåòèè ñíñòòíøò á ñíçäåíéè è ííäéòééåöèè íáúåéòíá.

Â íá-åéå ðåááíòû

Ííüçíååòåëü ñàí çåäååò íåòåäåííûå, òí áñòü äáííûå í ñòðóéòóðå áóäóùèõ íáúåéòíá. Ýòí ýäéýåòñý áéäååíííñòüþ ñèñòåìû: íðååééä ðåçìåòèè çåäåòþòñý ñàíèí ííüçíååòåëåí, à íá æåñòåí ñçíèñáíû á ïðíäðåííá. íåòåäåííûå ñíñòíýò èç ðåáéíííá è íáäñòðíáé íáá íéé. Øàáéíí (ééé èéèññ) ííæíí ííéíàòü èåé åáñòðåéòíûé òèí äáííûõ, ííðååééýþùèé áéä íáúåéòà. Íáíðéíáð, á ñòáíäåðòíûõ òåéñòåò, ñ êòòðûíé ðåáíòåò ñèñòåìà ðåçìåòèè, íðåäííéåååòñý òåééå ðåáéííû, éàé «ñòðåíéöà», «ñòðíéà», «ñéé íáíðòéå, éííéðåòíúá ñòðåíéöà, ñòðíéà èéé ñéíáíòðíå á òåéñòå – ýòí íáúåéòû ñííòåååòñòåóþùèõ ðåáéíííá.

Êàæäííó ðåáéííó

Ííèíèñàí ííðååééåííûé íááíð ííéé è íáðåíé-åíéé. Ñ íííñòüþ ííéé áíéé íáúåéòû á åééóíåíóå ííáóó áûou ñâýçáíû ñ äðóåéé. Óàé ííæíí ííèñàòü, ÷òí ñòðíéà òåéñòåò íòíñòéñý è éåéé-òí ñòðåíéòå, ñéíâá ðåñíííéååíû á ííðååééåííû ñòðíéåò, à áñýéåý ñéíâíòðíá èíååò ÷-åñòü ðå-÷. Ííéý ðåáééíá – ýòí íááíð òéííá íðéçíåééíá, éíòíðûá ííáóó áûou ó íáúåéòå ýòí íáí ðåáéííá. Íá íéó ííáóó áûou íáéíæåíû áñòåñòåååííûå íáðåíé-åíéý. Ýòé íáðåíé-åíéý íòíñòéñý è é òèíó äáííûõ çíà-åíéé ííéé, è é ñàíèí çíà-åíéý ííéé è éò ííáííéé (íáíðéíáð, åñéé ííáééåæåúåá á ðåáéííá «äéååíû -éáíû» – íóäåéüíá ñéíáíòðíá, èíåþùàý íáááæ, òí ýòí ìáááæ áíéæåí áûou èíáíéòåéüíû). Óàééá íáðåíé-åíéý çáíèñûåþòñý á áéäåå ííé-åñééò ñéííáéé íá ííéý (é èò ííáííéý ñ éþáûí óðíáíáí áéíæåíííñòé). N-éòååòñý, ÷òí áéý áñýééíá íáúåéòå äáíííáí ðåáéííá ýòé íáðåíé-åíéý äíéæíû áûou ðåéäååñòååíííéñòéíí.

Øàáéííû ííáóó

âûñòðåéååòñý á èåðåðòèè íáñëäååíáéé. Ýòá áíçííæíñòü íéàçûåååòñý í-åííû óáíííé íðé ííèñàíéé íåòåäåííûõ. Øàáéíí-íáñëååíéé íðéíáðåòååò áñå ñâíéñòåà (ííéý è íáðåíé-åíéý) ðåáéíí-íðååéé, áíáàåééyy è íéí, áíçííæíí, ñâíé íááíð ííéé è íáðåíé-åíéé. Øàáéíí-íðååéé ííæåò áûou ááñòðåéòíû (ò.å. èñïíëüçíååòñý á èá-åñòåå áíáùååí íðååéå äðóåéé ðåáéíí-íáñëååíéé). Ñíçäåååò íáúåéòû ááñòðåéòííí ðåáéííá íáéüçý. Íáíðéíáð, åñéé ííéüçíååòåëü ðí-åòí íááééòü áñå íáúåéòû ñèíòåéñè-åñééé ðåçìåòèè ííéáí «éíííáíòåéé», íí ííæåò ííðååééòü ýòí ííéáí ó íáúååíí ááñòðåéòííí ðåáéííá «ñèíòåéñè-åñééé íáúåéò» è áûâñòé èç ýòíáí ðåáéííá äðóåéé ðåáéííû.

Iàäñòðíéèà iàïïíèíååò
àáñòðàéòíûé ñàáéíí. Ííà ñòðíèòñý iàä óæå ñóùñòâóþùèìè ñàáéííàìè èéé
iàäñòðíéèàìè, éîòñòúå íàçùâàþòñý èàíæäàòàìè ià âõíæäåíèå à yòó iàäñòðíéêó.
Êàæäåíòò êàíæäàòò ñìçåòò áûòü ïðëíèñàíí ñòðíèåé ià ååñí âõíæäåíèå à iàäñòðíéêó.
Êàé e iãðåíè-åíèå ñàáéííà, yòí óñëíâéå iðåäñòåäéÿåò ñíáíé ëíàè-åñêíå áûòðàæåíèå,
çàâéñÿùåå ìò êííèðåòííà íáúåéòà, ååñ íïéåé, íïäííèåé è ò. ä. lïæíí
èíáóéòéàíí ïðåäåäéèòü ííýòéå ðåàéèçàöèè êííèðåòíù íáúåéòíí iàäñòðíéêé
èéé ñàáéííà. Åí-íåðåñò, åñÿééè íáúåéò ï ðåàéèçóåò ñíáíé ñíáñòååííûé
øàáéíí è åñå ñàáéííû-ïðåäéèé yòíåíí ñàáéíí. Äàéåå, tóñòü È – êàíæäàòò iàäñòðíéêé È iàúåéò ï ðåàéèçóåò È. Òíäàà ñ÷è
äéÿ íáúåéòà O áûññéíåíí óñëíâéå ià âõíæäåíèå
êàíæäàòà È à í.

lääänöödîéêà, êäè è
øàáëíí, iïæåò áûòù òèíí iïëý. Nïïòâåòñòâåííí, íáúåèò iïæåò áûòù cíà÷åíèäì òàéïäí iïëý, åñëe iï ðåàéëçóåò òàéïé òèí. Ýòî äàåò iïðåäåéåííóþ ãèäéîñöü à iïðåäåéåíèè iåðåäåííúö. Äí-åòïðûö, à iïðåäåííà èíååòñý âïçíïæíñöü iïðåäåéåíèè, ðåàéëçóåò èëé åäííûé íáúåèò óéàçàííóþ iåèñöödîééó, áûâåñòè ñïèñïé iåèñöödîåê, ðåàéëçóåíûö åäííû íáúåèòí, à òàéæå áûâåñòè åñå íáúåèòû, ðåàéëçóþùèå åäííóþ iåèñöödîééó. Nïàè ýòè íáúåèòû iïåóò èíåòû ðåçíûá øàáéííû; èò íáúåäéíýåò èëøü òî, ÷òî iïðè áûïïéíåèè ñïèñïé åòïæåäåíéy iú ïòïñèl èô ë èäíííé iåèñöödîéêå. Iïýòïíó íáèñöödîééè óäíáíí ðåññìàòðéåàòù èàé iïèñàíéy iïðñöûñ çàïðïñïâ è åäííû íáúåèòû, òî åñòù òàéëöö çàïðïñïâ, èïòïðûñ áîçåðåùàþò iòååéüíûé ñïèñïé íáúåèòû. Iïðè íåðì iïæåò ñïóæéöö íáèñöödîéêå «iïäéäæåùå», à èïòïðöþ åòïæéö «ñïèñïåòïðià» (iïðè íåëè÷ëé ñïèñïé íå iååäåæ), à òàéæå øàáéíí «iïëý» áåç ñïèñïé.

Âñýèéè íáúâéò èíáåò
íáÿçàòåëüíóþ òâéñòíâáóþ èííííáíó «ñíäåðæàíèå». Ñíäåðæàíèå íáúâéòà ííæåò
çääåàòùñý ííeüçíàòåëéàí, èéáí áû-ñéñëöùñý íí ïíðääåëéííù íðåâéèàí ÷åðåç
ñíäåðæàíèå ííéé. íáúâéòò èíåþò òàéæå ñíäöéèüñá äåññéðéíòðû äéý
ñíðòðéðíâéè è ñðåâíâíèå. ííé, á ÷-ñòíñòè, ííçâíéýþò ñòðíèòü çàíðñû è
íäðàíè-âíèå íá ííðÿâíè ñéíà (íàíðéíâð, íáéòè áñâ ñâýçè «ñóáñòáíòéâ-àòðèáóò», á
éíòíñòû ñóáñòáíòéâ íáòíäèòñý ðåíüøà áòðèáóò).

Ííéy ðàáéííá ííáóó
áúòù ððåôô áééäíâ: íáû÷íúâ ííéy, éíééâéööè è áéàïàçííú. ííéâ
éíééâéööè íðéé-àâðñý ìò íáû÷ííá ííéy ðâí, ÷òí íðâäííéâââðò ñðàçó
íâñéíéüéñ ðâçéé-íûð çíà-âíéé (íàïðéïåð, íâííðíâíû ÷ëáíû). Áéàïàçíí – yòí
«ñâýçíáy» éíééâéööè, ðí âñóü ííâæññôâ íáúâéòâ, éâðùéö ííáðyâ â ñíùñéâ
óííðyâ-âíéy íí áâñðéëïòíðâ. Áéy áéàïàçííâ áññòàòí-íí çâääâòù íá-àëüíûé è
éííá-íûé íáúâéò. Òëíè-íûé íðéïåð áéàïàçííá – ñðòðíè è ñòðâíéöâ èéè éâéèå-ëéâí
âññôâññôâíû ñâýçíûâ áíéüøéâ ððââñâíû òâéñòâ. ííéy ðàáéííá òâéæâ åääéyöñy íá íáýçàòâéüíûâ
è ííðéííáéüíûâ. íáýçàòâéüíá ííéâ çàííéyâðñý íðé ñíçääâíè íáúâéòâ
(íàïðéïåð, íðé ñíèíðâéñ-âññéé ðâçïåðéâ). Áéy ííðéííáéüíû ëíéâ íðâäéâââðñý
ñíèññéé âíçííæíûâ ââðéâíòâ çàííéíâíéy. Áâíûé ñíèññéé ðíðíèðâðñý íá íñííââ
íâðââé-âíéé ðàáéííâ è óâéâ çâííéíâíû ëíéâ.

Åñèè åñâå êåíäèåäåòù
ìåññöðîéèè èåläþò íåúèå ïïëý, ðî ïðè çàïèñè óñëïåêë ý à ïïëå òèïå ýòîé íåäñòðîéè
òåéñòåðå ïïëý ñåññòåðå òåéñòåðå ïïëý, ñåññòåðå òåéñòåðå ïïëý. Èñåñòå ðåéñòåðå ïïëý

ññâè (ñöèññäëüñâ) ññëÿ. Íaúåêò ïðèñáðåòàåò òàéñâ ññëå ñòðëüñî â òì ññö÷àå, åññëè íí ðääëëçóåò íaëññòðíééò. Ñ ñññùþ òàéñâ íaðåáíèçì ññëíí õäíñâ ññëñûñàòöù ññðòëññëè ññëóþ ðàçìåòëó (ýòî ïðèñáíÿëññü â áàçâ ãàíñû «Íaâññññññëàÿ íaðåáäÿ ëæòññü»).

Óñéíâéý è íâðáìè÷-âíéý
â íâðáàäííûõ çâääþòñý íà ñíâöèàëüíï ÿçûêå, êíðóðûne èíðåðíðåðòåðñý
íðíâðáìílï. lïëüçíàðåðäü ïíæåò êâé ñíçääâåàòü èô ñ ñíññüþ êíñòðóðòíðà
íðåðáíè÷-âíéé, òâé è çâëëñûâåàòü åðó-íþ. Bçûê ñíâðåðæòò íññíâíûå èíâð-âññéèå
íðåðåòòðû AND, OR,
NOT, íâðåðòòðû ðàâåâíñòâà (=), íâðåââíñòâà (<>) è
ñòðåâíâíèý, íðèíâëåðæíñòè (IN) è íáïðéíâëåðæíñòè
ííæåñòâó (NotIN). Å áûñðàæáíèý ñíñòð ó-âñòâíâåòü ïíéý
è èô ñíññüþ ñ éþáûõ óðíñíâí ãéïæáíñíñòè. lïéå-êíèéâéðòý âññâäà ðàññìàòðèåàåðñý
êâé ííæåñòâí, êðñíà õíñí, ííæåñòâí ïíæåò áûòü çâääíï ýâíï ñ ñíññüþ ôðæåðíûõ
ññéíâí è íâðå-èñëéâíèâí åôïíÿüèô ã íâðâí íáúâéòâí. Åò íðèíâð íñðâíè÷-âíéý íà
øðåâí «ñâýçû ñ ñíææññâííù àòðèåðòò»:

([Àòðèáóò].[xàñòü
ðå÷è] IN {"ïðèëåàòåëüíâ', "ïðè÷àñòèå'})

OR

((Àòðèáóò].[xàñòü
ðå÷è] = 'ìåñòîèìåíèå')

AND

([Àòðèáóò].[Ëèöî] NotIN {'1-å', '2-å', '3-å'})

AND

([Àòðèáóò].[Ëåêñåìà] NotIN {'è'}))

OP

([Àòðèáóò].[xàñòü ðå÷è] = '÷èñëèòåëüííå').

(Cääñü
òàéë çäìëñäííà óñëëâèå íà ëëöî àòðëáóòà íðñòî íçíà÷àåò, ÷òî ýòî ëëöî
íòñóòñòåáóò.)

Nðåäè îñíâíûõ

Ílāðáòìðáòìðá áýcúáéå áñòöù òàéæåá Þílåðáòìðá ïðílåâðéè ðåâæéçåöèè íáúâðòíí íàäñòðíééè
ééè ðááéëííà IS é óñeïâíúé Þílåðáòìðá IF (äey íáðåùåíéy è ítëyì íáúâðòíâ, êtòïðûå, áñíáùå áñâîñðy, íá
ýæéþþoñy íáúèìè). Ílñòü, íaïðèìåð,
ííéå «ílæéåæåùå»
ííéåð áûòü áûðåæåíí êæé ñeïlåîòìðíí, òæé è íóéåí. Ílñòü ñøáéëííù «ñeïlåîòìðíà» è «ííéü», áåçóñéïâíí, áññäyò á
íåéîòìðóþ íáññòðíééò. Ó ñøáéëííà «ííéü» íaò ítëyì «ílæéåæ»; ñílòâðòñòðåíí, ííéå
«ílæéåæ» áñòöù ñíëüêí ó ñøáéëííà «ñeïlåîòìðíà». Ílýòñò öñeïâéå íá «ílæéåæåùå»
ííæíí çàïèñåòöù òæé:

1F

([läääæàùåå] IS Ñëîâîôïðìà, [läääæàùåå].[läääæ]=[èìáíèòåëüíûé]).

Èíóâðôåéñ íðîâðàìù î ïðåäññòàâëåí
ïàíåëÿè ïáúâèòîâ. Ýòè ïàíåè è áûâàþò ðàçíûõ âèäåíâ; èõ îñííâíàÿ çàëà÷à –
ïòíâðàæåòû ñïåöèåëüíû ïáðàçì î ïðåäññòàâíû ïáúâèòû. Íà ãäííû ëìâîò ì
ïðîâðàìíà ïðåäññòàâíû òàêèå àêäû ïàíåëåé ïáúâèòîâ, êâè ïàíåëü ïàâèñàòè, ïàíåëü
îñííâíâíà ñâåñòà, ïàíåëü-ñïëñîè, ïàíåëü ñâîéñòà è ó. ä. Åçàèïñâÿçè ïàæäó íèè è èõ îñâåäåíèå ëìñûâàåòñÿ à
ïòíâðåëüíî xml-ôàééå.

Äëy ðàáìòû îïëüçîàòåëü
ïïæåò áùåäåëýòü á ìáíåëýö áðóïïû íáúåêòïâ. Èàæäàÿ áðóïïà íáíçíà÷àåòñÿ ñâïè
öååòï.

Ñèñòåìà áóäåò
ñíâåðøåñòåñòåíàòùñÿ äàëüøå. Íøåäññèàåòñÿ áíåäðèòü â íåå ìåðàíèciù

Ííèòåâòìàòè÷åñêîé ðàçìåòèè, ðàáîòû ñ «ïðåäëîæåéÿìè», ñïñíáû âèçóàëèçàöèè
ðàçìåòèè è ò. ä.

Ñièñîê
ëèòåðàòóðû

Çiáíèí è äð. 2006 – Çiáíèí,
Á. È. Óíèååðñàëüíàÿ ñèñòåíà ðàçìåòèè òåêñòà ATE-2
/ Á. È. Çiáíèí, Á. Á. Íàðêåëíà // Nîñðåííûå èíòîðìàöèííûå òåðííëåèè è
ïèñüíàííà íàñëåäèå: ìò äðåâíèõ ðóéíèñåé ê ýéåéòðííûì òåêñòàí : íàðåðèàëû
íàæäóíàð. íàó÷í. êííô., Èæåâñê, 13-17 èþëÿ 2006 á. – Èæåâñê, 2006. – Ñ. 51–55.

ObjectATE, a universal system for
text markup

Alexey I. Zobnin, Alexandra V.
Markelova

Vinogradov Institute of the Russian
Language of the Russian Academy of Sciences, Lomonosov Moscow State University,
Moscow, Russia

The object model and new features of ObjectATE,
a universal text annotation system, are described. This system allows a user to
define his own annotation models by describing classes, add-ins, fields, and
relations in the metadata layer.