

Èäáíẽíáẽý è òáõíéèá ðàçíáòèè á Òáõííéíáèè ñíáøáíííáí íááíðá

Ááõíð Áéáèñáíáð Áééòíðíáè- Èíááéáíéí
27.06.2008 á.
Íñéááíáá íáííáéáíéá 16.07.2008 á.

Òáçèñú á òíðíàòá RTF (57.69 kB 2008-07-13 21:38:16) Òáçèñú á òíðíàòá PDF (122.14 kB 2008-07-13 21:37:33)

Òáõííéíáẽý

ñíáøáíííáí íááíðá (ÒÑÍ) áñòú ñííñíá ðááíðò ñ òáèñòáìè, ñíááðæáúèèè òðááíáíóú ðàçèè-íúò ìèñúíáííúò ñèñòáì. Ñ òáõíé-áñéíé òí-èè çðáíéý ìíæíí áíáíðèòú ìðíñòí í òðááíáíòáð, òðááóòòèð ðàçííé èíòáðíðáòáòèè, ÷òí ñðàçó ááéááò òíðíóèèðíáéó ìðíáéáíú áíéáá íáúáé è áéòóáéúííé è áéý òáèñòáí, íáííðíáíúò ìí èñííéúçóáííé ìèñúíáííííòè. Áéý ðááíðò ñ òáèèè òáèñòáìè òðááóáòòèý ðáøáíéá òð, ò áçàèííñáýçáííúò çááá-: 1) ñíçááíéá ááéíé ìéàò-òíðíú áéý ìðááñòááéáíéý á íáíí òáééá ðàçííáí ðíáá òðááíáíòá; 2) íáòíæááíéá ñííñíáá ðàçáéáíéý òáèñòá íá òðááíáíòú; 3) ìðááíñòááéáíéá ìíéúçíááòáèò ñííñíáá ìíèñáíéý èíòáðíðáòáòèè òðááíáíòá.

Òáèü ááííé

ðááíðò – ñíñðááíðí-èòùñý íá áííðíñáð ðàçíáòèè è á, èíòáðíðáòáòèè, ðàññíòðáá èò íáúèá ñííááíéý, è ìíéàçáòú áíçííæíííòè áéý ìíáñáíááíé ðááíðò èññéááíáòáéý ìðéíýòíáí á ÒÑÍ ìíáòíáá á ñðááíáíé ñ áðóáèè èìáòòèèèñý ìíáòíááìè è ðàçíáòèá, ìðáæáá áñááí XML.

1. Òáðíèíú.

Òáèñò ìú ìíéíááí èáè çàòèèñèðíááííóò ìíñéááíáá-òáéúííòú çíáéíá. Òáèñò ìíááðáááòòèý èíòáðíðáòáòèè ÷áéíááéí èèè ááòííàòí ñ èííéðáòíé ìðáèòè-áñéíé èèè èññéááíáòáéúííéíé òáèüò. Èáæáíá èññéááíááíéá ìíèðááòòèý íá èíòíðíáòèò, ñáýçáííóò ñ èííéðáòíúèè ìáñòáìè á òáèñòá, è áñýéáý ñíáá ñííñíáá èíòáðíðáòáòèè íóæáááò-ñý á óéàçáíéè ìáñòá òáéíé ñíáíú. Ñí-áí-éòí-ííñòú óéàçáíéè íá ìáñòá á òáèñòá, ñáýçáííúò ñ èííéðáòíé èíòáðíðáòáòèéá, ìú íàçúáááí ðàçíáòèíé. Òáèè ìáðàçíí, ñ òáèñòí ìíæáò áúòú ñáýçáíí ìííáí ðàçíáòíé, ìðáá-áòòèè ðàçíúí òáéýí.

Ááæíúèè

ìðéíáðáìè ðàçíáòèè ýáéýòòèý ðàçáéáíéá òáèñòá íá ñíñòááéý-òúèá ááí òðááíáíòú è ìðááéáíéá ñðóòèòòú òáèñòá á áéáá èáðáðòèè. Ñáýçú èáðáðòèè ñ òáèñòí ìáòñéíáéáíá áçáéýáíí íá òáèñò (òáèüò èíòáðíðáòáòèè), òáè ÷òí ðàçíúá ðàçíáòèè ìáííáí òáèñòá ìíáòó ìíááðæéááòú ðàçíúá èáðáðòèè. Íí á íáúáí ñéó-áá ðàçíáòèá – ýòí íá ðàçáéáíéá è íá èáðáðòèý, á òíéúéí íááíð óéàçáíéè íá ìáñòá á òáèñòá.

Íí ñííñíáó ñáýçè ñ òáèñòí ðàçíáòèè ááéýòòèý íá ááá èèáññá: áíáøíéá è áíóòðáíéá.

Áíáøíýý

ðàçíáòèá – ýòí íááíð ñíñúéíè íá ýéáíáíòú òáèñòá, òðáíéíúé ìòááéúíí ìò òáèñòá. Òáèáý ðàçíáòèá ìáèéó-òèè ìáðàçíí ñííòááòòèòáòáò íáçááèñèíííòè íáúáèòá (òáèñòá) ìò ááí ìíèñáíéý, íááñíá-èááý íáçááè-ñèíííòú ðàçíáòíé. Ìðéíáðíí áíáøíáé ðàçíáòèè ìíæáò áúòú ìáðá-áíú ìáñò á ìíñéááíáá-òáéúíííòè çíáéíá, áúðáæáííúé ìíáðáìè ñéíá òáèñòá, áóéá ñéíáá.

Áí-íáðáúó,

éíóáðíðáòáòèy òáèñòá ííèðá-áòñy íá òíèúéí íá yáíóþ ðàçíáòéó ñòáíáàðòííáí òíðíàòá, íí è íá íáyáíóþ, èáé è íá èpáúá áíííéíèòáèúíúá ñèíáíèú èèè èííáéíáòèè ñèíáí-èíá. Íá íóæááþòñy á ñíáòèáèúíúó íáòéáò è íñíáúá ñèó-áè éíóáðíðáòáòèè, èááíòèòèòèðóáíúá íí ñí-áòáíèþ íáòéè ñ yéáíáíòàìè èñóíáííáí òáèñòá, òáé èáé òáèííó ñí-áòáíèþ ííæíí íðèáàòú íñíáóþ éíóáðíðáòáòèþ òáí æá íáðáíèçìí, -òí è ñáííé íáòéá.

Áí-áòíðúó,

òíòy yéáíáíóú ðàçíáòéè íèèðúááþò ñíáñòááííúá òðááíáíóú òáèñòá, è íèí íá íðááúyáèyáòñy ñòðíáíá òðááíáíéá éíáòú íáðííé (“çáéðúááþúéé”) yéáíáíò (-òí íá íáòááò èñíèú-çí-áàòú íáðíííòú íðè íáíáóíáèííòè, ííááðæéáày á, íðè ííèñáíèè èíóáðíðáòáòèè).

Yòè

ðáðáíèy ááèáþò ðàçíáòéè á òáèñòá òáíðáòè-áñèè íèíèíáèúíúíè: íáòéá ñòááèòñy òáí è òíèúéí òáí, ááá áíçíéèááò íáíáóíáèííòú á ñíáíá éíóáðíðá-òáòèè, íá óñíàòðèáááíáy èç áðóáèò íáðèáðíá. Á ðáçóèúòáòá íáèíí-ðúá çááá-è ííáóò ðáðáòúñy ááç áíáñáíèy áíííéíèòáèúííé ðàçíáòéè á óæá èíáþúéáñy éíðíóñá.

Á. Ííæáñòááíííòú

ðàçíáòéí íááñíá-èáááòñy áíçííæíííòúþ íáðáá èàæ-áúí ííá-áííúí òðááíáíòíí áñòááèòú ííáúé, ñèóæááíúé òðááíáíò ñí ñáíáé íáò-éíé. Òáèéá òðááíáíòú ñèíáíèñè-áñèè íá íòèè-áþòñy íò òðááíáíòíá èñóíáííáí òáèñòá (è ííáóò íáñòè èpáíá éííè-áñòáí éíòíðíáòèè, ñíááðæáòáèúííé äèy ííááðæéááííáí ááííé ðàçíáòéíé áñíáèòá). Íðè yòíí, áñèè íáòéá òðááíáíòá íá íðèááí íèèáéíáí çíá-áíèy íðè ííèñáíèè èíóáðíðáòáòèè, òí ááñú òðááíáíò (á íá òíèúéí íáòéá) èáííðèðóáòñy. Òáèèí íáðáçìí, áíáááèáíéá á òáèñò ííáúò íáòíè íá áèèyáò íá óæá ñóúáñòáóþúéá éíóáðíðáòáòèè, á ííèñáíèè èíòíðúó íðí yòè íáòéè íè-ááí íá ñèáçáíí.

Á. Íáçááèñèííòú

ðàçíáòéí ðááèèçóáòñy, èðííá óèáçáííé ñííñíáíííòè íáðáíèçíá éíóáð-íðá-òáòèè èáííðèðíáòú “-óæéá” òðááíáíòú, òáí, -òí íí (íáðáíèçí) áíííòèááò è íðíèçáíèúííá çáááááíéá òðáèòíáèè ñáííáí ðáçáéáíèy íá òðááíáí-òú, ííçáíèyñ ñ-èòáòú ñóúáñòááíííúè äèy ðáçáéáíèy òíèúéí íóæíúá íáòéè.

Á. Íðíñòíòá

ííèñáíèy éíóáðíðáòáòèè ðàçíáòéè íááñíá-èáááòñy òáí, -òí ñ òí-èè çðáíèy ðáááè-òí-ðá éíðíóñá ííá ñíñòíèò á áúíèñúááíèè á íðíñòíí òáèñòí-áíí òáééá (“ñòèèááíí òáééá”) äèy èáæáíé íáòéè ñííñòááèáíèy áé òáíí-èè ñèíáíèíá è íðáíáðá-çí-ááíèé, èíòíðúá ñèááóáò íðèíáíèòú è òðááíáíò. Íðáíáðáçíáíèy, í èíòí-ðúó èá, ò ðá-ú, óèáçúááþòñy ñáíèí èíáíáí, çá èíòíðúí ííæáò ñòíyòú èéáí íðíáðáííáy óóíèòèy, èéáí íáíííáðáíííá íðáíáðáçíáíèéá, çáááááííá íáðá-íáí íðíñòúò èèè ðáéòðñèáíúó çáíáí á ñíáòèáèúííé òáèñòíáíé òíðíá. Ñéíóáèñè-áñèè èíáíá òáò èèè áðóáèò íðáíáðáçíáíèé íá ðáçèè-áþòñy, òáé -òí ðáááèòíð éíðíóñá ííæáò íí-ñáíáíò íáðáííðáááèyòú èáèèá-òí èç ñèñòáííúò óóíèòèè.

Íðèíáðàìè

íáíðĩáðàìííúò ìðáíáðàçĩááíéè ìĩáóò ñéóæèòú: ìáðááĩá ìðááñòàáéáíèý ñèìáíèĩá á óó èèè èíòò ððèòòíáòò èíáèðĩáéò, óááèáíèá óááðáíéè, ðáññòàíĩáèá ìáðáíĩĩá è/èèè ììðòáííúò áðáíèò, èáììàòèçàòèý, ñĩýòèá/ðáññòàíĩáèá òèòè è áðóáéá ìĩáðàòèè, ñĩáòèòè-íúá äèý ìðòíáðàòè-áñéíé òðááèòèè è ìðááñòàáéýòèá ñèíæíĩñòú íá ñòíèüèí àèáíðèòìè-áñéòò, ñèíèüèí ìðááìáòíí-ñĩáòèàèüíòò. Íáñĩòðý íá òáóíè-áñéòò ìðĩñòíòò, ñàìè ìðáíáðàçĩááíèý, ìèò-àáíúá òàèèì íáðàçĩ, ìĩáóò áúòú áíáíèüíí ñèíæíúìè, áí áñýèì ñéò-àá çà íèìè íáú-íí ñòíèò áíèüøíé òðóá ñĩáòèàèèñòà, èíòíðúé ìðááñòàáéýáò ñĩáíé çàéíí-áííúé ðàçóèüòàò. Íðááèàáááìàý òíðìà àà, ò áíçĩíæíĩñòú ñĩáòèàèèñòàì íáíáíèáàòúñý è òàèèèè ðàçóèüòàòàìè.

3. Íðĩòáññ

èíòáðĩðáòàòèè ðàçĩáòèè. Ñòèèááíé òáéé ìðèìáíýáòñý íá áòíðì ÿòàĩá èíòáðĩðáòàòèè. Íáðáúé è òðàòèé ÿòàì çááàòò òñèíáèý è òáèü ðááíòú ðàçðááíò-èèá ðàçĩáòèè.

Ìáòèè

ñòáíáàðòíĩáí àèáà íá-èíáòòñý ñĩáòèàèüíúì ìáðéáðì “íá-àéí” (-àñòì ÿòì çĩáè “%”). Íáðáúé ÿòàì èíòáðĩðáòàòèè ñòááèò ìáðáá íèìè áííèíèòàèüíúé ìáðéáð “èííáò”, òàè ÷òí íá áðáíèòàò òðááíáíòíá òáçúáááòñý èííòáìèíáòèý ÿòèò ìáðéáðíá. Íá òðáòúáì ÿòàíá óááéýòòñý áñá ó-áñòèè òáèñòà ìáæáò ÿòíé èííòáìèíáòèáé è áèèæàèøèì è íáé ìáðéáðì “èííáò”. Ííýòíò çááà-a áòíðĩáí ÿòàìá, èíòíðúé ììðááéýáò ñĩñòàáèòáèü ñòèèááíĩáí òáéèá, – ìðĩèñàòú á í, òàèèá çàìáíú, á ðàçóèüòàòà èíòíðúò íóæíúá òðááìáíòú íá òíèüèí íá óááèèèèñú áú íá òðáòúáì ÿòàìá, íí è ñĩááðæàèè áú á ñááá ìðááíèñáíèý íóæíúò ìðáíáðàçĩááíéè. Èñĩèüçóý á òàèèò çàìáíáò è ñàìè ìáðéáðì “íá-àéí” è “èííáò”, ìĩæíí áèáéí áèèýòú ìá ñèíæèáðááñý á ðàçóèüòàòà ìáðáñá-áíèý ðàçĩáòíè ìáðáíá-àèüíĩá “ðáçáéáíèá”.

ÒÑÍ

ðááèèçĩááíá íá ìèàòòíðìá PHP, ÷òí áúèí áúçááí á, ìáðáíá-àèüíúì ìáçĩá-áíèáì – ìðááñòàáéýòú á èíòáðĩáòà á ðàçĩúò áèáàò áðááíáðòññèèá ìáñĩĩáíèý. Ñĩáòèòèèáòèè, íáíáçĩáèìúá äèý ìáðáíĩá Óáòíèíáèè ìá áðóáéá ìèàòòíðì, ìðááñòàáèáíú ìá ñáèòà ìðĩáèòà “Óĩá çĩáìáíúò ìáñĩĩáíèé” (URL: <http://znamen.ru/tsn/tsnp-zam.php>).

Ideology and technique for encoding using the mixed scripts technique

A. P. Ershov Institute of Informatics Systems of the Siberian Branch of
the Russian Academy of Science, Novosibirsk, Russia

This paper discusses common problems in text encoding and proposes their solution using a “mixed scripts technique” (MST), a way of working with texts containing parts which need different interpretation. MST’s way of text markup, together with a procedure of markup interpretation, is given in comparison with other solutions like XML. A special form of non-programmed user-defined transformations is shown, which can become common for carrying out certain kinds of scientific research.