

# ONE WORLD, ONE WIKI!

C. Scott Ananian <cananian@wikimedia.org> [[User:cscott]] (Wikimedia Foundation)

## Language Converter

MediaWiki uses [[mw:LanguageConverter]] to automatically transliterate articles between closely related languages or dialects or script variants of a language or dialect. It is used on 11 wikis, and has been requested on about 35 more. Here are some examples of conversion pairs:

<div><div><div><div><div><span></span></div><div>English</div></div><div><div><span><span></span></span></div><div>(American/British)</div></div></div><div><div><div><span></span></div><div>LanguageConverter not used.</div></div><div>Spelling and usage differences exist between American English, British English, Indian English, and others.</div></div></div></div>	<div><div><div><div><span></span></div><div>[[Elevator]]</div></div></div></div> <div><div><div><span></span></div><div>An <b>elevator</b> is a type of vertical transportation that moves people or goods between floors (levels, decks) of a building, vessel, or other structure.</div></div></div>
---	--

## Native Variant Editing

LanguageConverter is oriented to readers: it converts the article text unidirectionally into readable text in a consistent variant. But as soon as a user begins to edit, they are confronted with the source text in a mix of variants, as illustrated by the intermingled Cyrillic and Latin scripts in the article from Serbian Wikipedia shown below. This mixture of scripts can be a huge barrier to editing in communities where individuals are typically only fluent in a single variant.

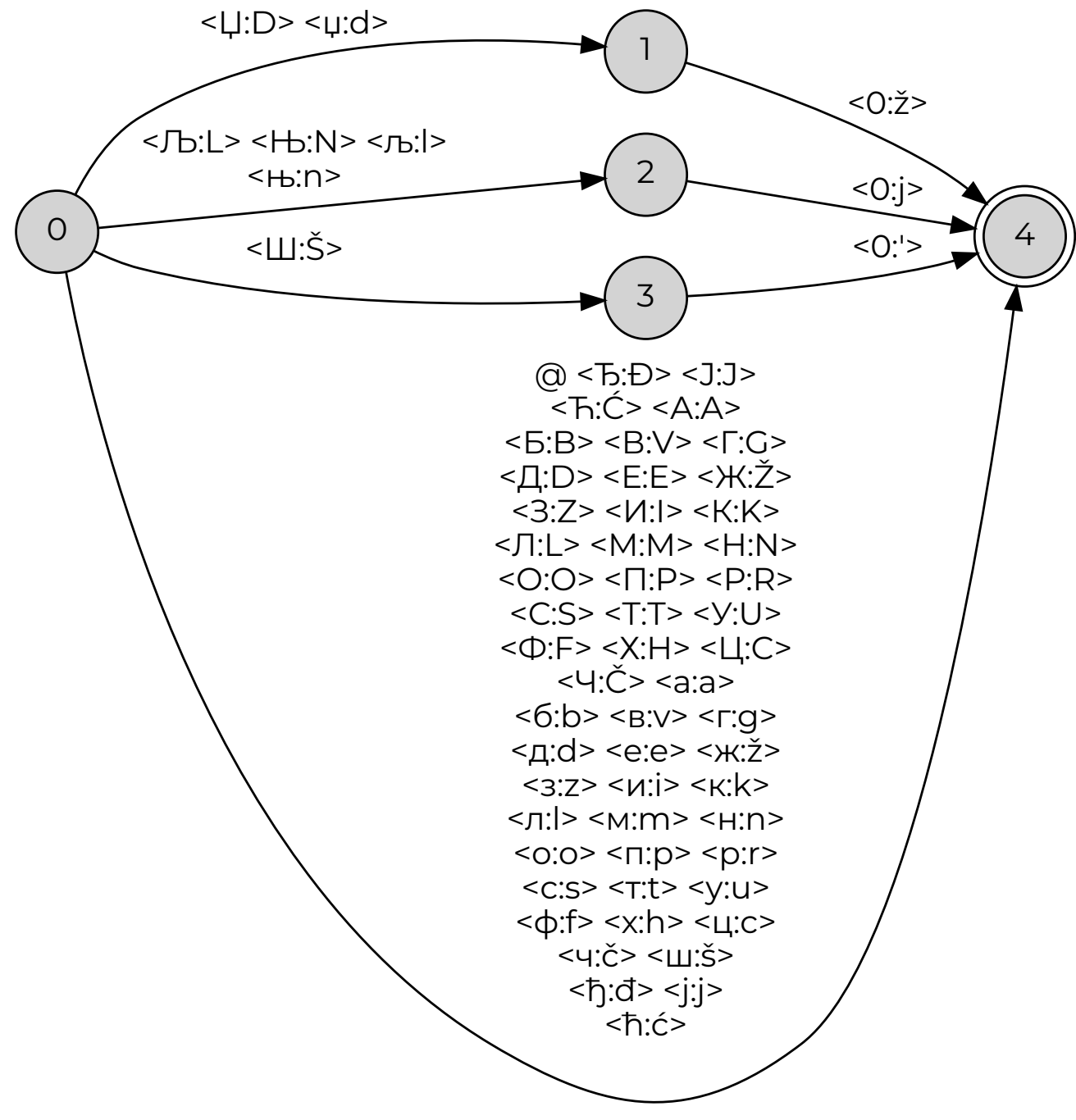
The Parsoid team has been experimenting with a new bidirectional implementation of LanguageConverter, based on Finite State Transducers (FSTs). These allow automatic annotation of wikitext such that it can be round-tripped to its original variant losslessly. With these annotations, an Wikimedian can edit an article in their preferred consistent variant.

Unedited portions of the article will round-trip to their original variant, preventing dirty diffs, and only edited sections will reflect the variant which the editor saw. On wikis where the community has chosen to author all articles in a single variant, all text can be losslessly saved as the chosen variant, regardless of which variant the editor used.

**We can make editing easier on wikis using LanguageConverter!**

Списак 118 познатих хемијских елемената
<span></span> <div>Следећа табела садржи 118 познатих хемијских елемената.</div> <div><ul style="list-style-type: none"><li><b>Атомски број, име, и симбол</b> служе независно као јединствени идентификатори.</li><li><b>Имена</b> су она која су прихваћена од стране <b>IUPAC</b>; провизиона имена за недавно произведене елементе који нису формално именовали су дата у заградама.</li><li><b>Група, периода, и блок</b> се односе на позицију елемента у <b>периодном систему</b>. Бројеви група су у тренутно званично прихваћеној нотацији; за старије алтернативне нотације погледајте <b>Група периодног система елемената</b>.</li><li><b>Stanje materije</b> (<i>Čvrsto, tečno, ili gasovito</i>) се односи на standardne uslove <b>temperature i pritiska (STP)</b>.</li><li><b>Pojavljivanje</b> прави разлику između elemenata koji се јављају у природи, kategorisane kao bilo <i>Praiskonski</i> ili <i>Prolazni</i> (у смислу raspada), i <i>Sintetički</i> elementi koji су произведени tehnološkim путем, i нису prirodno poznati.</li><li><b>Opis</b> sumira својства elementa користећи опширне kategorije које су присутне у periodnom sistemu: <b>aktinoid, alkalni metal, zemnoalkalni metal, halogen, lantanoid, metal, metaloid, plemeniti gas, nemetal, i prelazni metal</b>.</li></ul></div>

Serbian Wikipedia article showing mixed Cyrillic/Latin script



Partial FST for Serbian Cyrillic-to-Latin conversion

## Translation Suggestion Tool

lala

## Zero-Shot Translation

lala