

# 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>English</div> <div>(American/British)</div> <div><i>LanguageConverter not used.</i></div> <div>Spelling and usage differences exist between American English, British English, Indian English, and others.</div>	<div>[[Elevator]]</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>[[Lift]]</div> <div>A <b>lift</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>Serbian (Latin/Cyrillic)</div> <div><i>LanguageConverter in use.</i></div> <div>Speakers are fully functionally digraphic, using both Cyrillic and Latin scripts. There are also vocabulary differences between Ikavian, Ekavian, and Ijekavian dialects which are not currently converted.</div>	<div>[[Лифт]]</div> <div><b>Лифт</b> је уређај за транспорт људи или терета међу спратовима зграда или радних платформи. Уобичајно се креће уз помоћ електромотора који или покреће узад за вућу и противтежни механизам, или пумпа хидрауличну течност за подизање цилиндричних клипова.</div>	<div>[[Лифт]]</div> <div><b>Лифт</b> је уређај за транспорт људи или терета међу спратовима зграда или радних платформи. Уобичајно се креће уз помоћ електромотора који или покреће ужад за вучу и противтежни механизам, или пумпа хидрауличну течност за подизање цилиндричних клипова.</div>
<div>Chinese</div> <div>(Simplified/Traditional)</div> <div><i>LanguageConverter in use.</i></div> <div>Simplified used in mainland China, Singapore, and Malaysia. Traditional used in Taiwan, Hong Kong, Macau, and among Overseas Chinese. Most speakers monographic; few can fluently proofread text in both variants.</div>	<div>[[电梯]]</div> <div><b>电梯</b>，亦称升降机、垂直电梯。在香港、新加坡和马来西亚俗称“𨋖”（英语lift的译音），是一种垂直运送行人或货物的运输工具。</div>	<div>[[電梯]]</div> <div><b>電梯</b>，亦稱升降機、垂直電梯。在香港、新加坡和馬來西亞俗稱「𨋖」（英語lift的譯音），是一種垂直運送行人或貨物的運輸工具。</div>
<div>Hindi/Urdu</div> <div><i>LanguageConverter not used.</i></div> <div>Urdu and Hindi are dialects of the Hindustani language, written in very different scripts: Arabic on the Pakistan side of the border, Devanagri on the India side. (Punjabi is a similar case, with four scripts used.)</div> <div>Currently separate small wikis; could combine efforts.</div>	<div>[[उत्थापक]]</div> <div><b>उत्थापक</b>, उच्चालितर अथवा एलिवेटर (lift या elevator) एक युक्ति है वस्तुओं एवं व्यक्तिओं को उर्ध्व दिशा में चढ़ाने-उतारने के काम आती है। प्रायः किसी बहुमंजिला ऊँचे भवन, जलपोत एवं अन्य संरचनाओं में उत्थापक लगा होता है जो गोलों को या सामान आदि को एक मंजिल से दूसरी मंजिल या एक स्तर से दूसरे स्तर पर लाता और ले जाता है। उत्थापक प्रायः विद्युत मोटर द्वारा चलते हैं।</div>	<div>[[رافع]]</div> <div>انتصابی نقل و حمل کی کل۔ جدید عمارتوں، ج۔ اڑوں اور کانوں میں استعمال۔ ونے والی تمام کھلی اور بند ساختوں اور لگانار چلنے والے ان پٹوں کو بھی رافع یا (Elevator:انگریزی) کی ا جاتا ہے جو بھاری چیزوں کو ایک جگہ سے دوسری جگہ پہنچاتے ہیں۔</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 this article from Serbian Wikipedia to the right. 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.

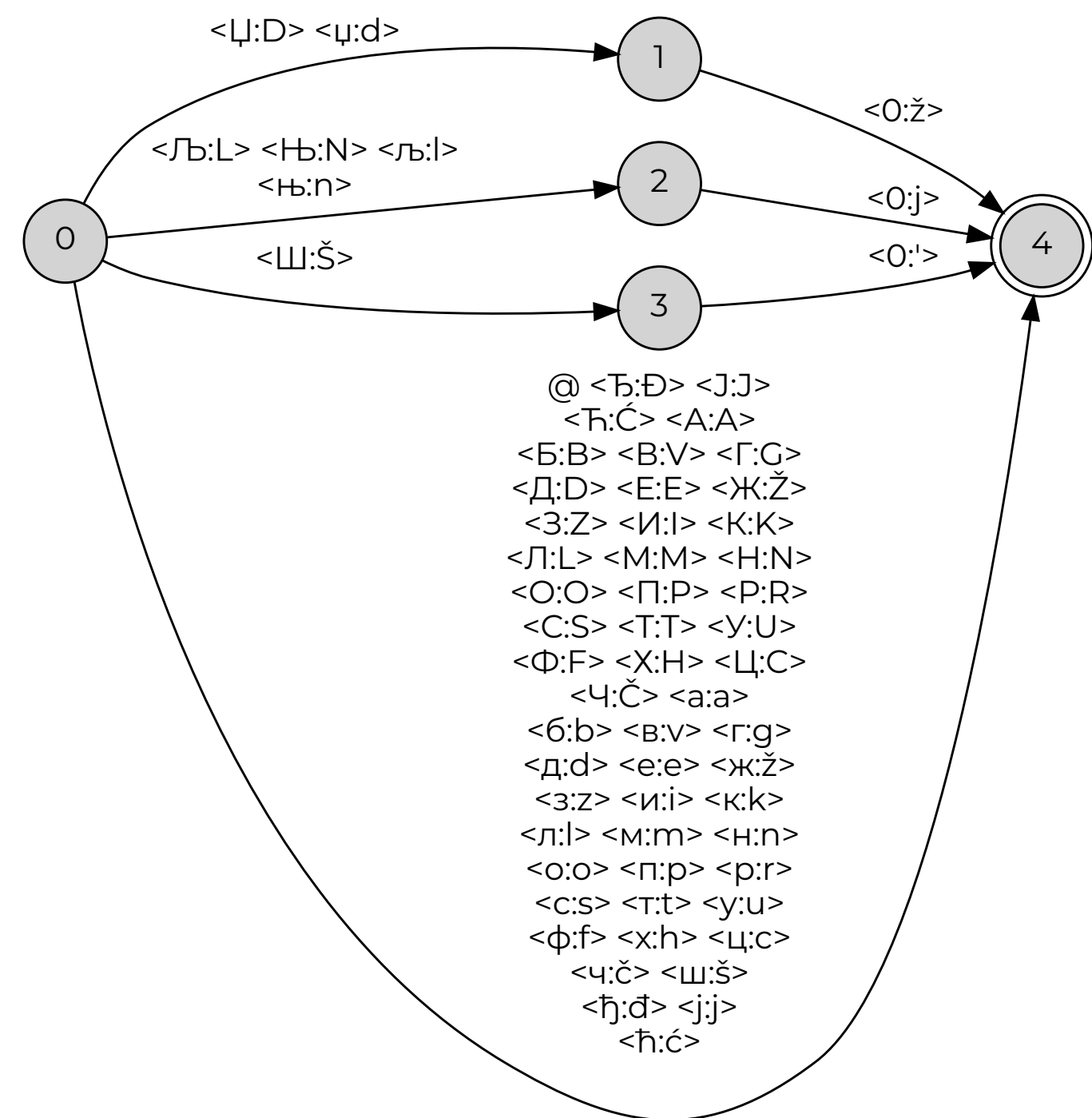


Fig. 1: Partial FST for Serbian Cyrillic-to-Latin conversion

## Translation Suggestion Tool

lala

## Zero-Shot Translation

lala