

# Editing challenges on multi-script wikis

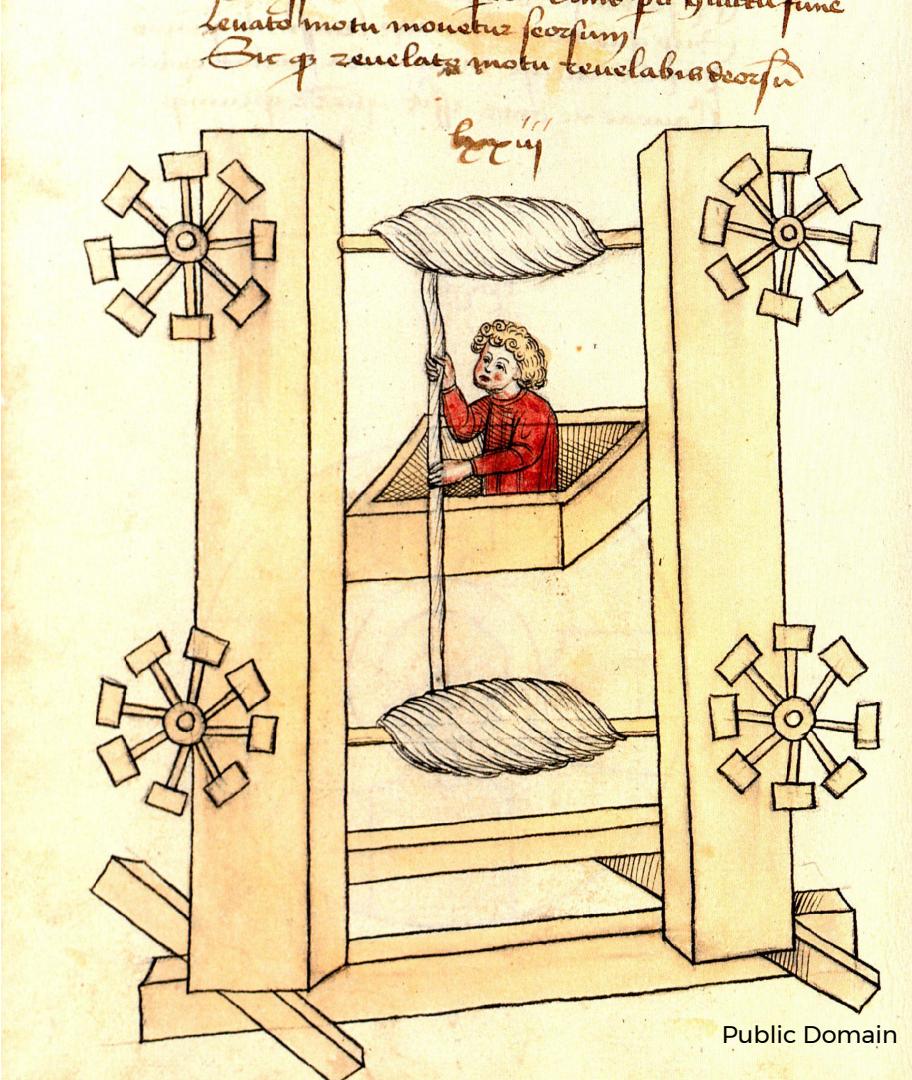
C. Scott Ananian  
Wikimania 2017



August 12, 2017  
Montréal, Québec, Canada

# [[Elevator]]

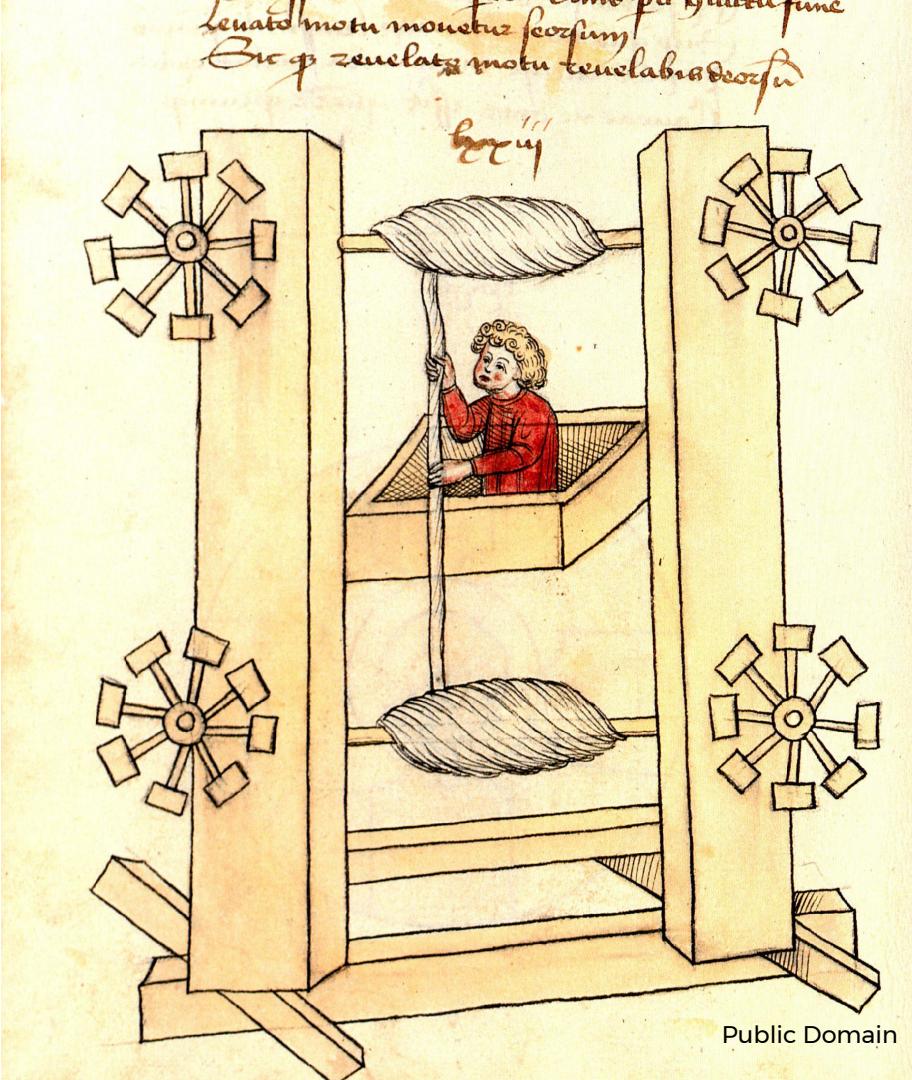
Elevator design  
by the German  
engineer Konrad  
Kyeser (1405)



Public Domain

# [[Lift]]

## Lift design by the German engineer Konrad Kyeser (1405)



Public Domain

# [[Lift]]

Lift je uređaj za transport ljudi ili tereta među spratovima zgrada ili radnih platformi.



WIKIMEDIA  
FOUNDATION



CC BY-SA 3.0 Peregrine981

# [[Лифт]]

Лифт је уређај за транспорт људи или терета међу спратовима зграда или радних платформи.



# [[電梯]]

電梯，亦稱升降機、垂直電梯。在馬來西亞、新加坡和香港俗稱「粒」(lift的譯音)，是一種垂直運送行人或貨物的運輸工具。



# [[电梯]]

电梯，亦称升降机、垂直电梯。在马来西亚、新加坡和香港俗称“粒”(lift的译音)，是一种垂直运送行人或货物的运输工具。



# [[उत्थापक]]

उत्थापक, उच्चालित्र अथवा एलिवेटर (lift या elevator) एक युक्ति है जो स्तुओं एवं व्यक्तिओं को ऊर्ध्व 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।



Public Domain

# [[رافعه]]

انتصابی نقل و حمل کی گل جدید عمارتوں، جہازوں اور کافتوں میں استعمال ہونے والی تمام کھلی اور بند ساختوں اور لگاتار چلنے والے ان پیٹوں کو بھی کہا جاتا ہے جو **Elevator**: رافع یا (انگریزی بهاری چیزوں کو ایک جگہ سے دوسری جگہ پہنچاتے ہیں طاقت سے چلنے والے رافع جو عام طور پر بھایاں سے کام کرتے تھے انیسویں صدی عیسوی سے استعمال ہو رہے تھے جبکہ اس صدی کے اوآخر مینبر قی رافعہ عام ہو گیا۔



Public Domain

**LanguageConverter  
converts words,  
scripts, and (if you  
stretch) languages**



# It brings our wikis together

(Although fighting over orthography is part of the fun?)



**Conversion  
in use on 11  
wikis;  
wanted on  
~35 more**

[\[\[meta:Wikipedias in multiple writing systems\]\]](#)

- Chinese
- Serbian
- Kazakh
- Kurdish
- Inuktitut
- Anglo-Saxon
- Shilha
- Tajik (partial)
- Uzbek (partial)
- Gan (partial)
- Cantonese (client-side)





But it can make editing  
harder.



CC by SA 4.0, C. Scott Ananian

# Mixed script editing

## Списак 118 познатих хемијских елемената [\[ уреди \]](#)

Следећа табела садржи 118 познатих хемијских елемената.

- **Атомски број, име, и симбол** служе независно као јединствени идентификатори.
- **Имена** су она која су прихваћена од стране **IUPAC**; провизиона имена за недавно произведене елементе који нису формално именовани су дата у заградама.
- **Група, периода, и блок** се односе на позицију елемента у **периодном систему**. Бројеви група су у тренутно званично прихваћеној нотацији; за старије алтернативне нотације погледајте [Група периодног система елемената](#).
- **Stanje materije** (*Čvrsto, tečno, ili gasovito*) se odnosi на standardne uslove **temperature i pritiska (STP)**.
- **Pojavljivanje** први разлику између елемената који се јављају у природи, категорисане као било *Praiskonski* или *Prolazni* (у смислу распада), и *Sintetički* елементи који су произведени технолошким путем, и нису природно познати.
- **Opis** сумира својства елемента користећи општине категорије које су prisutne u periodnom sistemu: **актинoid, alkalni metal, zemnoalkalni metal, halogen, lantanoid, metal, metaloid, plemeniti gas, nemetal, i prelazni metal**.

# Status



# Parsoid support

# Visual Editor support



CC by SA 4.0, C. Scott Ananian

# Reading variants

- Google indexing
- Kiwix offline reader
- Android app
- PDF
- REST API



# Future



# Source language annotation

Current implementation relies on character set to implicitly mark source language.

For reliable operation in some variants, we need explicit hints.

Ideally we can use an [annotation service](#), and not direct wikitext markup.

Initial work will focus on wikis where articles are stored in a consistent variant; for example, articles on yue.wiki stored in yue-hant variant.

Лифт => Lift

I think I have learned that the best way to lift one's self up is to help someone else.

— [Booker T. Washington](#)

# Single-variant editing

Parsoid/Visual Editor use [\*\*Selective  
Serialization\*\*](#) in order to convert only the edited portion of a page back to wikitext.

Same technology can be used to edit in your native variant but avoid disrupting the variant used in unedited portions when saving changes.

Please lift those boxes and place them in the elevator.

Please -{lift}- those boxes and place them in the lift.

# Glossaries

Word-based conversion dictionaries can get large and accumulate topic-specific entries.

Chinese wikipedia uses [templates](#) and a [custom gadget](#) to add new conversion rules.

Improve support by bringing this into core as article [glossaries](#).

**[[Glossary:Harry Potter]]**

-{en-uk:Sorceror;  
 en-us:Philosopher}-  
-{en-uk:Minister for Magic;  
 en-us:Minister of Magic}-  
-{en-uk:sherbet lemon;  
 en-us:lemon drop}-

# Content Translation

LanguageConverter can be viewed as a simple implementation of machine translation.

At the moment we only use machine translation to generate the initial revision of articles.

A future [Content Translation Tool](#) could be used to keep parallel texts in different languages (or variants) in sync.



# TL; DR



CC by SA 4.0, C. Scott Ananian

# Summary

- Not everyone uses Latin letters!
- Some texts which look different are really quite similar.
- LanguageConverter is a tool for sharing content when texts are closely related.
- **We can make editing easier on wikis using LanguageConverter!**
- These wikis could be good incubators for more general cross-wiki machine translation tools.



# THANK YOU

[\[\[commons:File:Editing\\_challenges\\_on\\_multi-script\\_wikis.pdf\]\]](#)



# The graveyard of unused slides



# Example markup

Marking up text which has variants:

- Bidirectional rules:  
-{zh-hans:computer; zh-hant:ELECTRONICBRAIN;}-
- Unidirectional rules:  
-{HUGEBLOCK=>zh-cn:macro;}-
- Disable conversion:  
-{R|SimpTrad}-

Also a bunch of stateful options for adding/removing rules (more on this later), and for working around title limitations.

See [https://www.mediawiki.org/wiki/Writing\\_systems/Syntax](https://www.mediawiki.org/wiki/Writing_systems/Syntax)

# Glossary Templates

From [\[\[zh:鋼鐵人3\]\]](#) ([\[\[en:Iron Man 3\]\]](#)):

```
 {{noteTA
|G1=Movie
|G2>Show
|G3=美国漫画
|1=zh-hans:罗伯特; zh-tw:勞勃; zh-hk:羅拔;
|2=zh-hans:奧德利奇·齊連安; zh-tw:奧德奇·齊禮安; zh-hk:奧
德奇·齊禮安;
|3=zh-hans:羅德斯; zh-tw:羅德; zh-hk:羅德;
}}
```

# Transliteration code

```
class SrConverter extends LanguageConverter {  
    public $mToLatin = array(  
        'а' => 'a', 'б' => 'b', 'в' => 'v', 'г' => 'g', 'д' => 'd',  
        'ђ' => 'đ', 'е' => 'e', 'ж' => 'ž', 'з' => 'z', 'и' => 'i',  
        [...]  
        'Х' => 'H', 'Ц' => 'C', 'Ч' => 'Č', 'Џ' => 'Dž', 'Ш' => 'Š',  
    );  
    public $mToCyrillics = array(  
        'а' => 'а', 'б' => 'б', 'с' => 'ц', 'č' => 'ч', 'ć' => 'ћ',  
        'д' => 'д', 'дž' => 'џ', 'đ' => 'ђ', 'е' => 'е', 'ф' => 'ф',  
        [...]  
        'Њ' => 'Њ', 'њ' => 'њ', 'Н!ј' => 'Нј', 'N!J' => 'НЈ'  
    );
```

# Word conversion code

```
$zh2Hant = array(  
    '侃' => '侃',  
    '𠂇' => '𠂇',  
    [...]  
    '; 克制' => '; 剋制',  
    '? 克制' => '? 剋制',  
);  
  
$zh2Hans = array(  
    '𠂇' => '倾',  
    '𠂇' => '𠂇',  
    [...]
```

