



*Enhancing diversity in the content creation process
at Wikipedia*

WP5

*Mathias Schindler (WIKI), Delia Rusu (JSI),
Fabian Flöck (KIT)*

Project review Y1

Luxembourg, 1st of December, 2011



1. Overview
2. Problem scenarios
3. Solution including reuse of R&D results
4. Mockup demo
5. Outlook



Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

Outlook

Overview



- The goal of Wikimedia's case study is to support Wikipedia editors in maintaining and improving the site, and to support readers in understanding the quality and biases of a given article.
- We are creating tools and extensions to support editors in the management, understanding, and decision-making about complex and heated controversies on Wikipedia.
- We want Wikipedia to offer high quality articles on both highly visible and as well as on more obscure topics.



- Not everyone contributes to Wikipedia
- Standard authorship has a strong spot at
 - Male
 - Academic background
 - 25-50 years of age
 - Developed countries
- Bias is a vicious circle
- Wikimedia Foundation targets for 2015 include:

„ Support healthy diversity in the editing community by doubling the percentage of female editors to 25 percent and increasing the percentage of Global South editors to 37 percent“



- Defining the scope of Wikipedia boils down to editorial decisions such as: An article for every episode and character and location of "The Simpsons" vs. a single article on the entire TV series.
- These decisions will influence:
 - the audience composition
 - the authorship recruitment
 - the internal and external perception of the project



Article [Discussion](#)

Read [Edit](#) [View history](#)

List of locations in *The Simpsons*

From Wikipedia, the free encyclopedia



It has been suggested that this article or section be [merged](#) into [Springfield \(The Simpsons\)](#). ([Discuss](#)) *Proposed since August 2011.*

This **list of locations in The Simpsons** includes towns, cities, businesses, [states](#), and other locations created for the television series *The Simpsons*. This list only includes places that exist solely in the *Simpsons* universe, not real-life locatic



- Variations in the number of „people of Muslim faith“

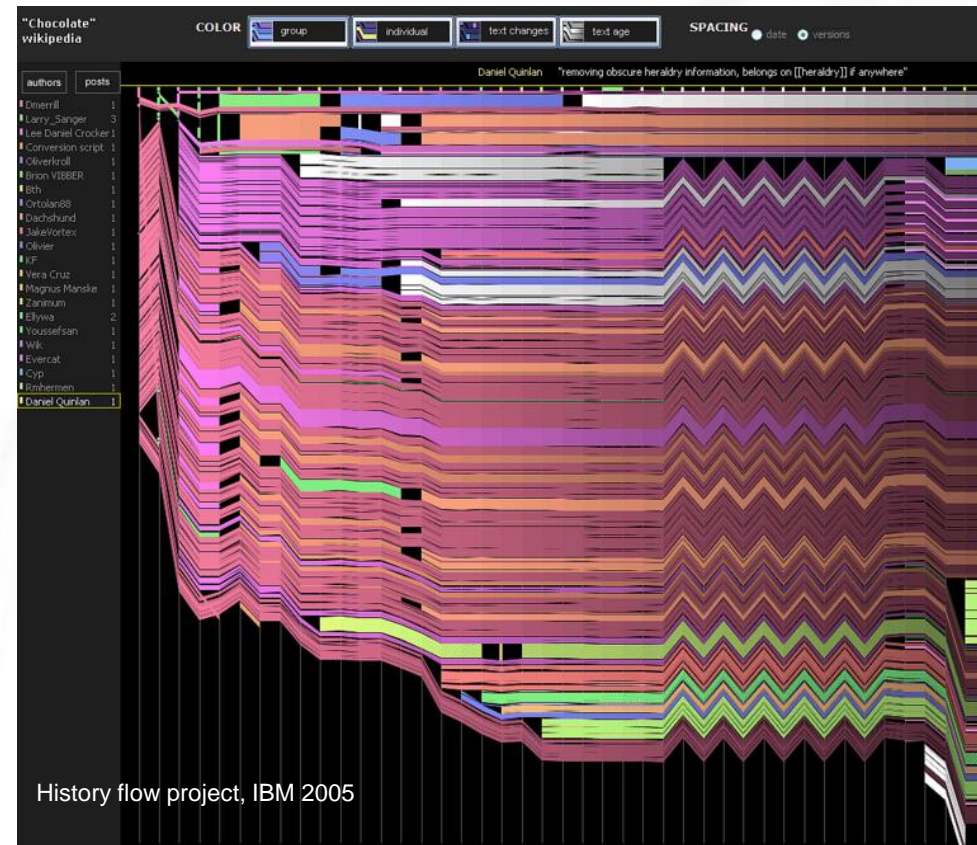
Lang	Lemma	Figure (in bn)
EN	Islam	Over 1.5
EN	Major religious groups	1.3 – 1.65
EN	Claims to be fastest-growing religion	1.57
HE	Islam	1.4
LB	Islam	1.1 – 1.5
ID	Islam	1.25 – 1.4

- Variations in borders drawn on maps in various Wikipedia language editions (Kashmir, Sea of Japan)

Diversity challenge (4) – editor behaviour



- Certain editing behaviours can lead to biased articles
 - e.g. a dominant editor group in an article that wins an edit war, 'pushing out' minority views
- Newcomers and 'outsiders' to an article can encounter problems adding content, especially mayor changes





- Definition of use case scenarios
- Collection of existing approaches in quantitative metrics for quality assessment
- Collection of bias-inducing editor behavior patterns and development of methods to detect them
- Metric definition in order to evaluate the development of Wikipedia article quality
- Engagement with the Wikipedian and Wikimedian community to explain the scope and the public benefit of the RENDER project
- Participation in Wikimedia community events to outline the current state of the RENDER project and to invite feedback from scientifically minded authors



Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

Outlook

Problem scenarios



Main Goal:

- Improvement of the quality, the value and the trustworthiness of Wikipedia by supporting Wikipedia users (readers and editors)

Use Case Scenarios:

- UC1: Display warnings to the reader when detecting bias
- UC2: Notify authors that an article needs to be updated
- UC3: Lower the barrier for readers to extend and/or correct articles

UC1: Bias detection and notification



- A regular visitor to Wikipedia opts into a tool that will display warnings whenever an article is shown with detected bias.
- The user is given a summary of the detected bias and detailed information on how the bias warning was triggered.
- The user is now in a position to engage in article editing to fix or amend the article in order to improve its quality and remove the biased parts.



- A retired professor of linguistics with advanced Wikipedia expertise and good standing as an author has committed herself into maintaining the entire topic in Wikipedia.
- Using a dedicated tool, she is given a list of articles that show signs of being outdated, incomplete or biased.
- The professor is now able to maximise impact of her work, focussing on the most deserving articles in her field of expertise.



- A student has a long time history of passively reading Wikipedia articles in the course of his studies.
- Dedicated tools are now providing him with information to understand which facts are missing in the article and are offering him resources that contain the missing information
- The user is now given a clear path to turn passive involvement into active and productive participation



Overview

Problem scenarios

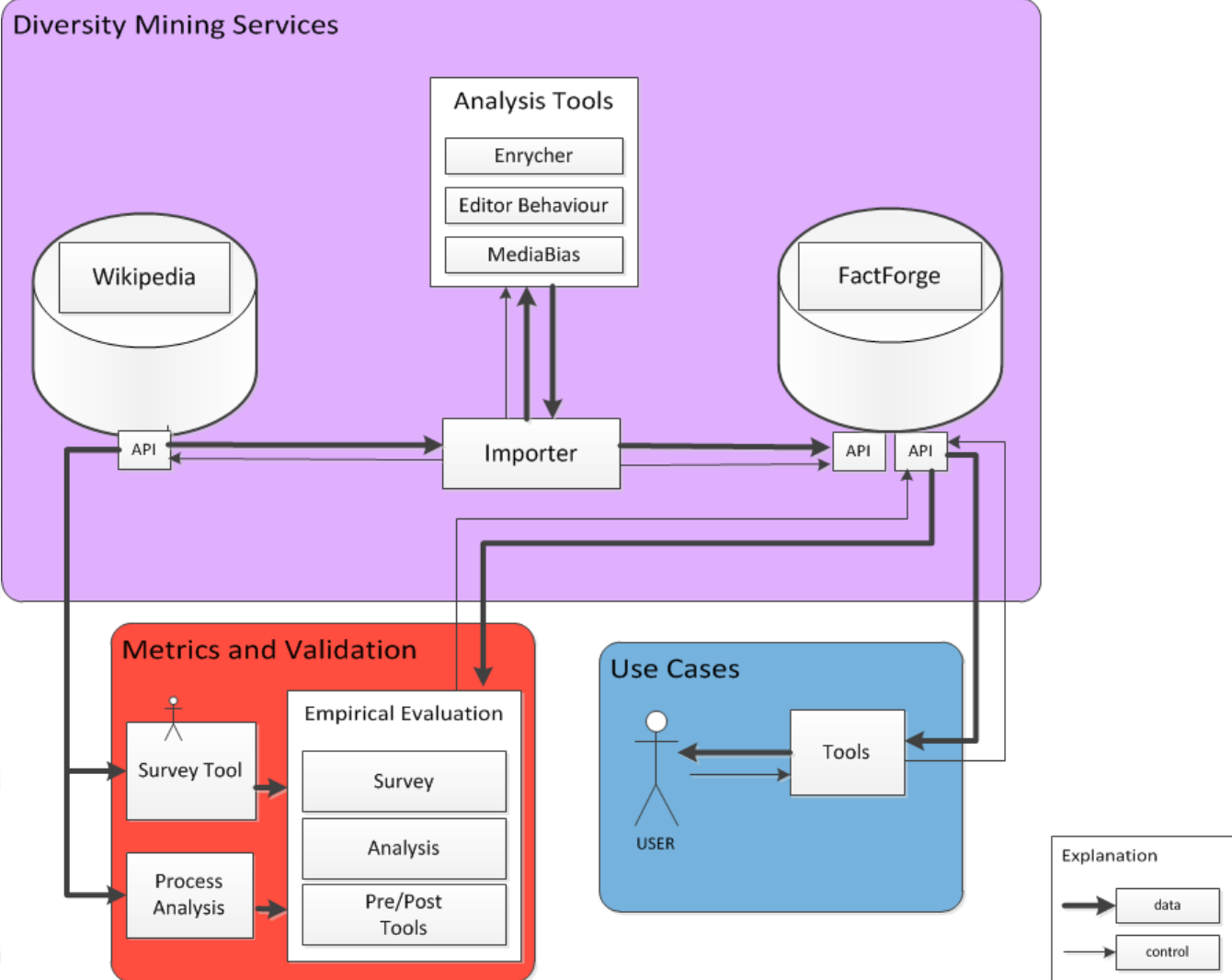
Solution including reuse of R&D results

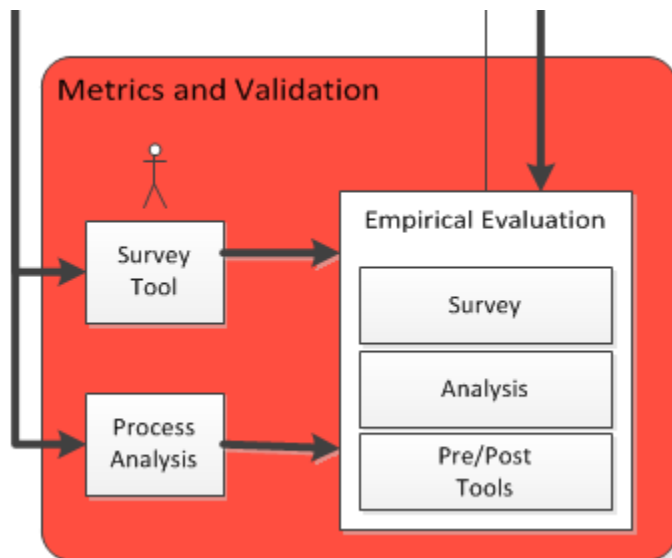
Mockup demo

Outlook

Overview of solutions including reuse of R&D results

Overview of solutions including reuse of R&D results





Evaluation of results from Enrycher, behavioral analysis and others

2 approaches to use Wikipedia's assessment expertise:

- assessment survey with Wikipedia users
- analysis of templates and the results of WMF Article Feedback Tool



- Analysis of Wikipedia's content development and quality
 - Fact coverage/ completeness:
 - Article length (number of words) compared to articles in other language versions
 - Number of articles which
 - ✓ have a bigger fact coverage compared to other language versions
 - ✓ have a lack of facts compared to at least one external source like a news article
 - Timeliness
 - Objectivity



- Analysis of Wikipedia's content development and quality
 - Fact coverage/ completeness
 - Timeliness:
 - Number of edits per day compared to the average
 - Number of articles which are
 - ✓ not reverted during the last day/week
 - ✓ without reverts but high editing in at least five language versions
 - ✓ out-dated compared to publishing dates of external sources (at least five days older)
 - Objectivity

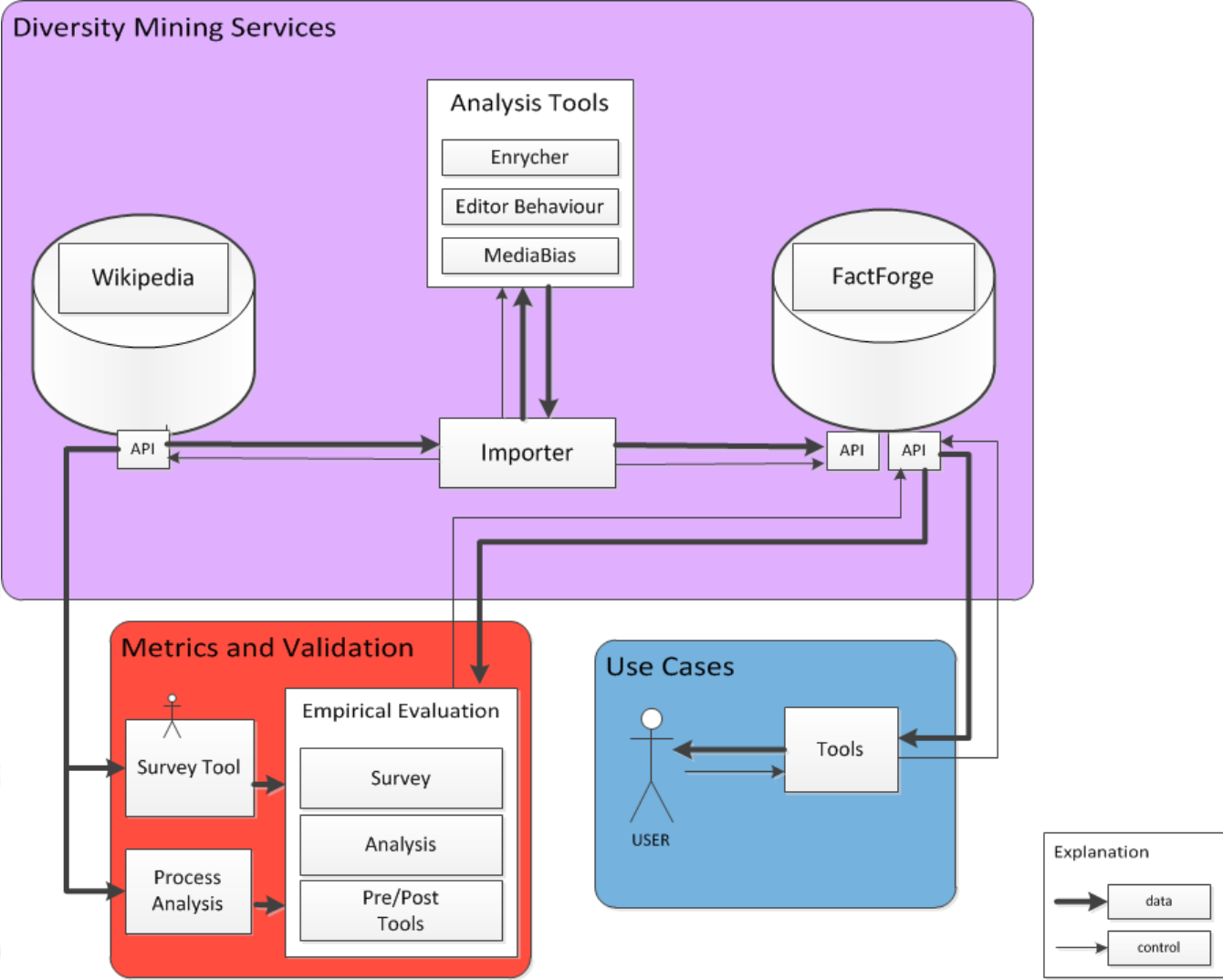


- Analysis of Wikipedia's content development and quality
 - Fact coverage/ completeness:
 - Timeliness
 - Objectivity:
 - Number of articles:
 - ✓ containing subjective words or expressions
 - ✓ identified as opinionated by JSI's algorithms
 - ✓ classified as opinionated by containing biased references



- Analysis of article-based editor behavior patterns and their development
 - Existence of editing patterns indicating bias:
 - Measured based on
 - ✓ Correlations of the chance of an edit getting reverted with editor, edit and article features
 - ✓ Social editor network metrics like centrality, clustering, density, etc.
 - ✓ Specific combination of behavioral mechanisms detected

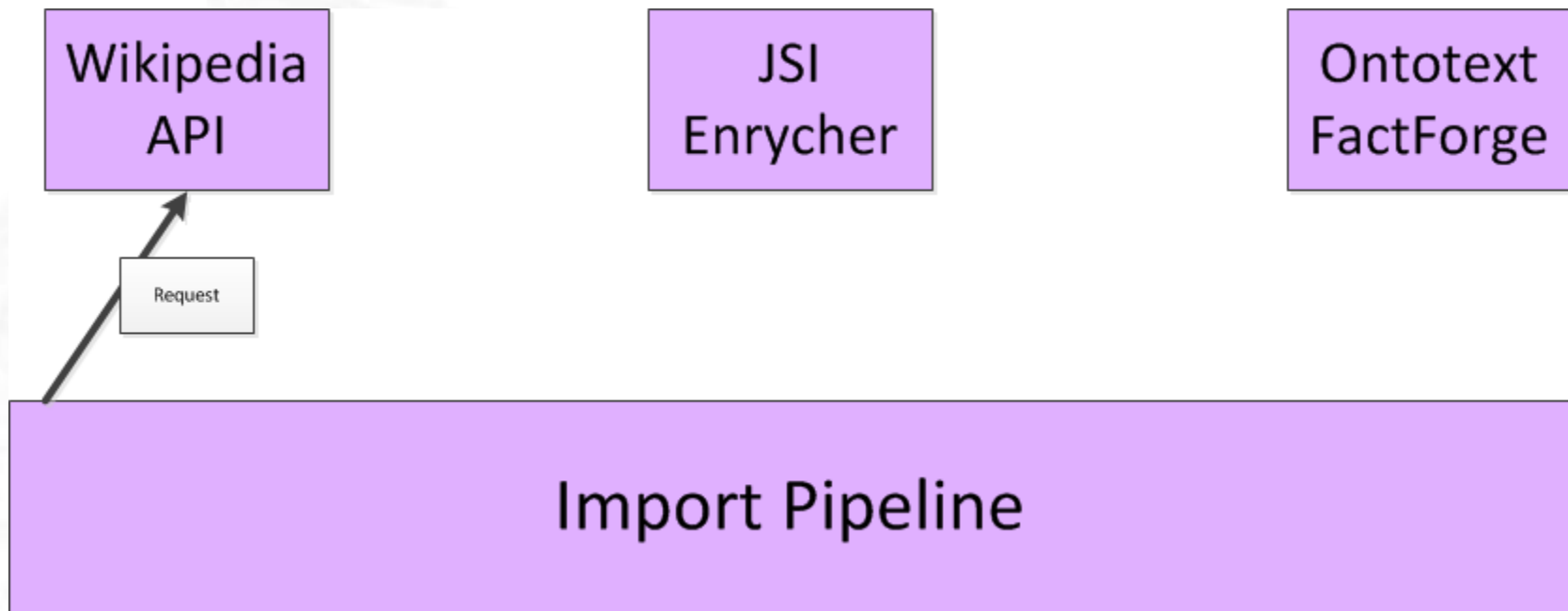
Overview of solutions including reuse of R&D results





Example for the import pipeline procedure:

API request: “Kalmar”

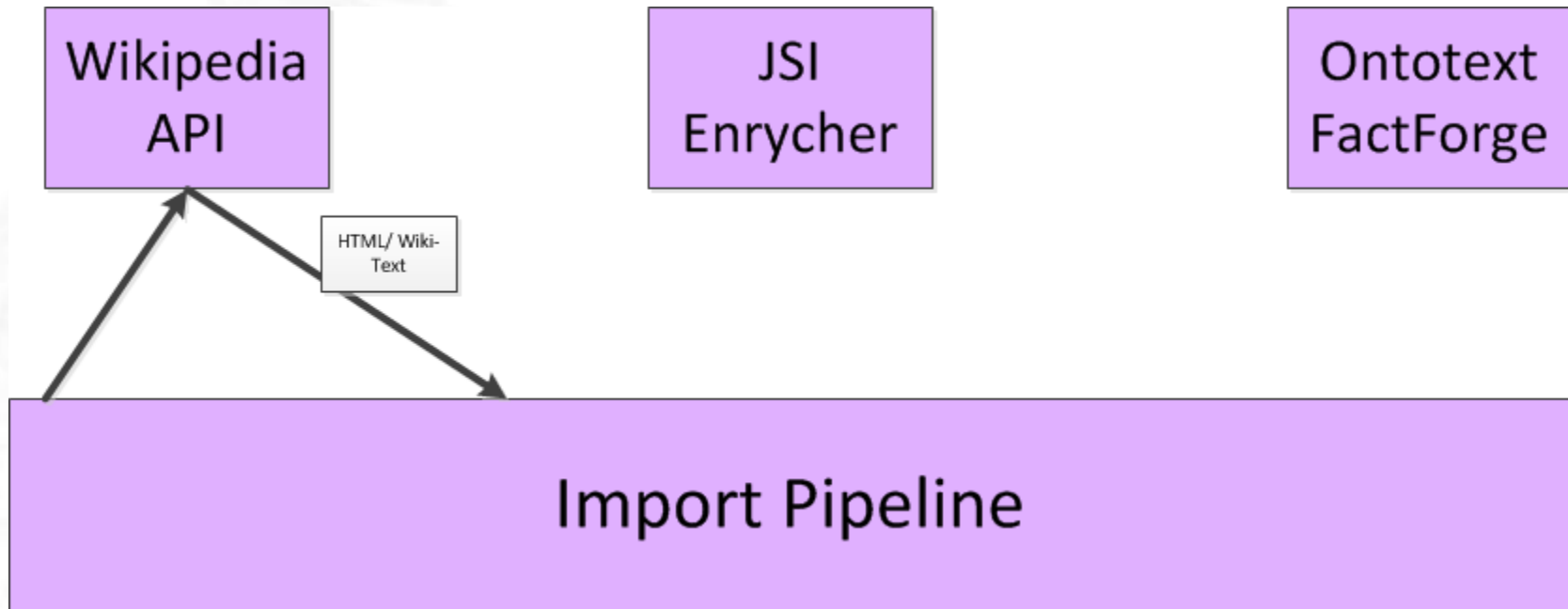


Diversity mining services – import pipeline



"Kalmar" is a [[cities of Sweden|city]] in [[Småland]] in the south-east of [[Sweden]], situated by the [[Baltic Sea]]. It had 62,767 inhabitants in 2010<ref name="scb" /> and is the seat of [[Kalmar Municipality]]. It is also the capital of [[Kalmar County]], which comprises 12 municipalities with a total of 233,776 inhabitants (2006).

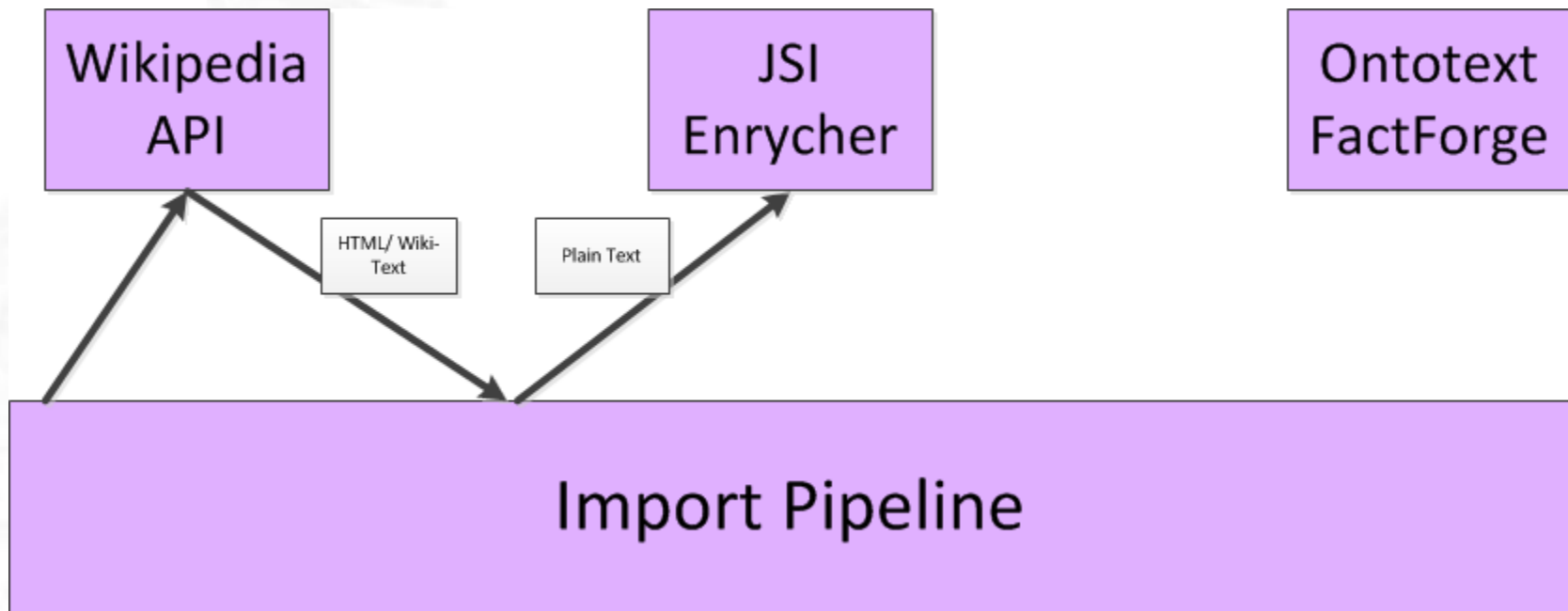
...





Kalmar is a city in Småland in the south-east of Sweden, situated by the Baltic Sea. It had 62,767 inhabitants in 2010 and is the seat of Kalmar Municipality. It is also the capital of Kalmar County, which comprises 12 municipalities with a total of 233,776 inhabitants (2006).

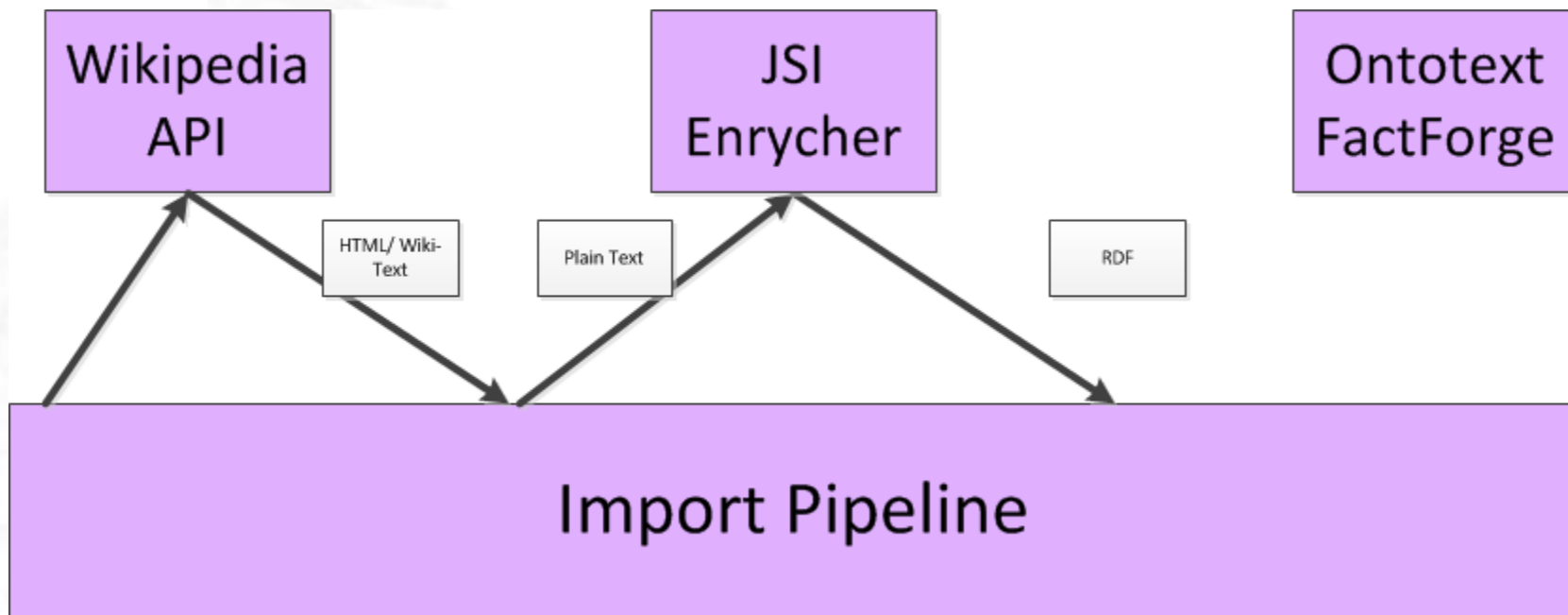
....



Diversity mining services – import pipeline



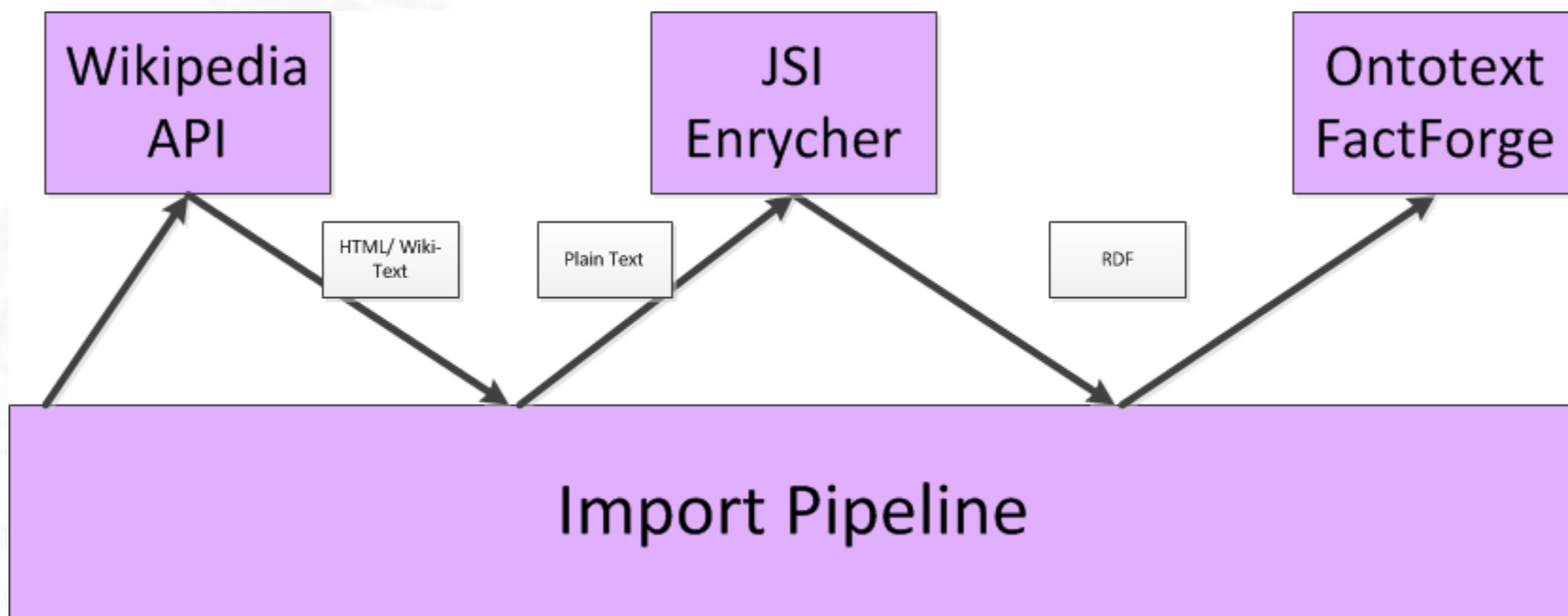
```
...  
<rdf:Description rdf:about="urn:document-3cbc5995-1679-4dad-9d2c-87bffb9bb69f">  
<rdf:type rdf:resource="http://xmlns.com/foaf/0.1/Document"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Regional/Europe/Sweden/Kalmar_Co  
nty/Localities/Kalmar"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Reference/Museums/Transportation/  
Maritime/Europe/Sweden"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Regional/Europe/Sweden/Maps_and  
_Views"/>  
...
```



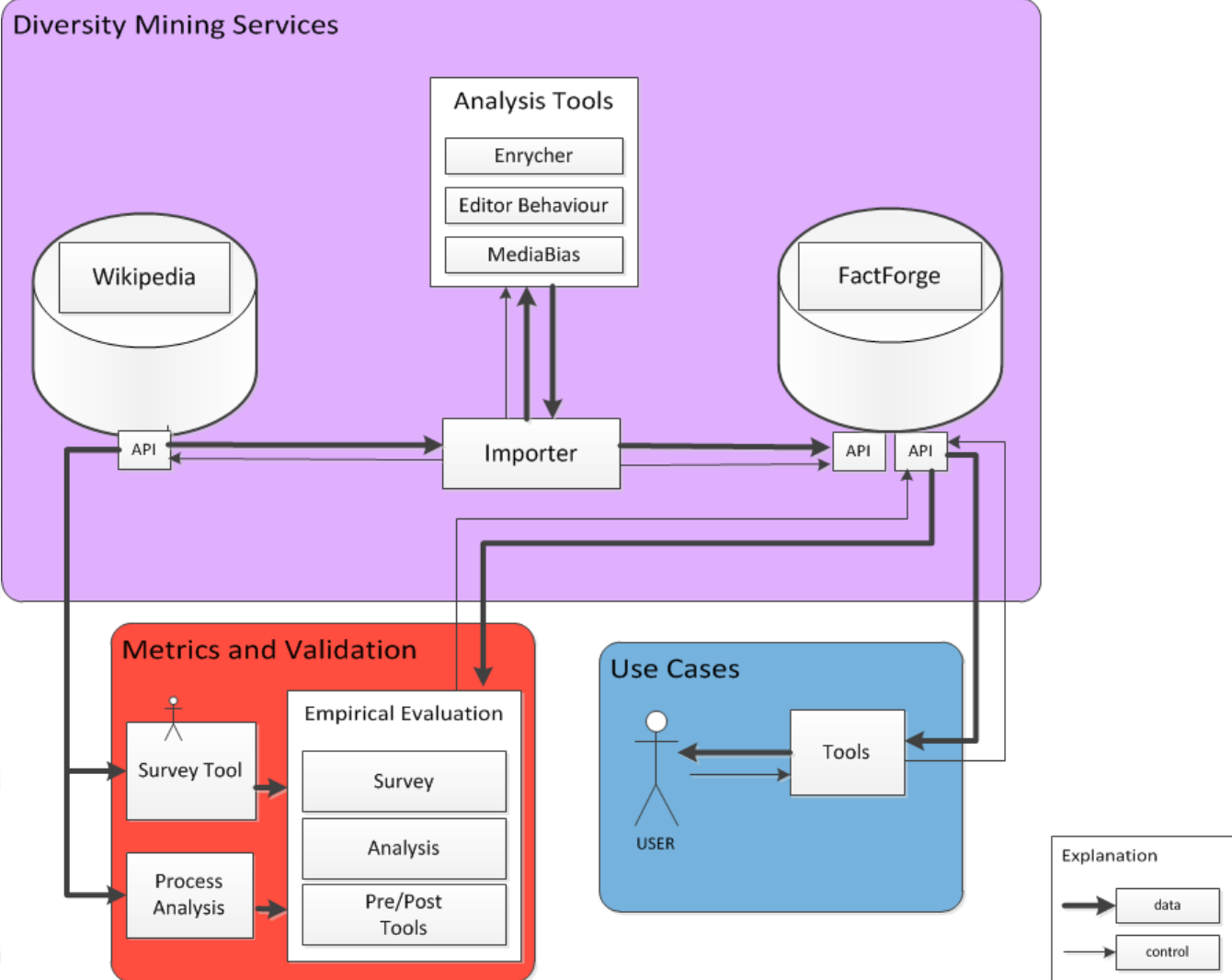
Diversity mining services – import pipeline



```
...  
<rdf:Description rdf:about="urn:document-3cbc5995-1679-4dad-9d2c-87bffb9bb69f">  
<rdf:type rdf:resource="http://xmlns.com/foaf/0.1/Document"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Regional/Europe/Sweden/Kalmar_Co  
nty/Localities/Kalmar"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Reference/Museums/Transportation/  
Maritime/Europe/Sweden"/>  
<dmoz:topic rdf:resource="http://www.dmoz.org/Top/Regional/Europe/Sweden/Maps_and  
_Views"/>  
...
```



Overview of solutions including reuse of R&D results





Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

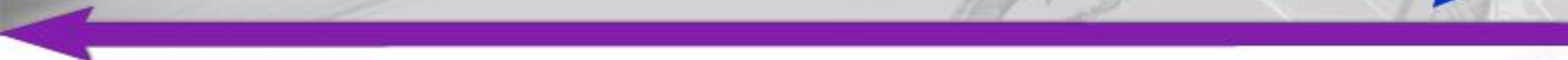
Outlook

Opinionated Wikipedia Articles

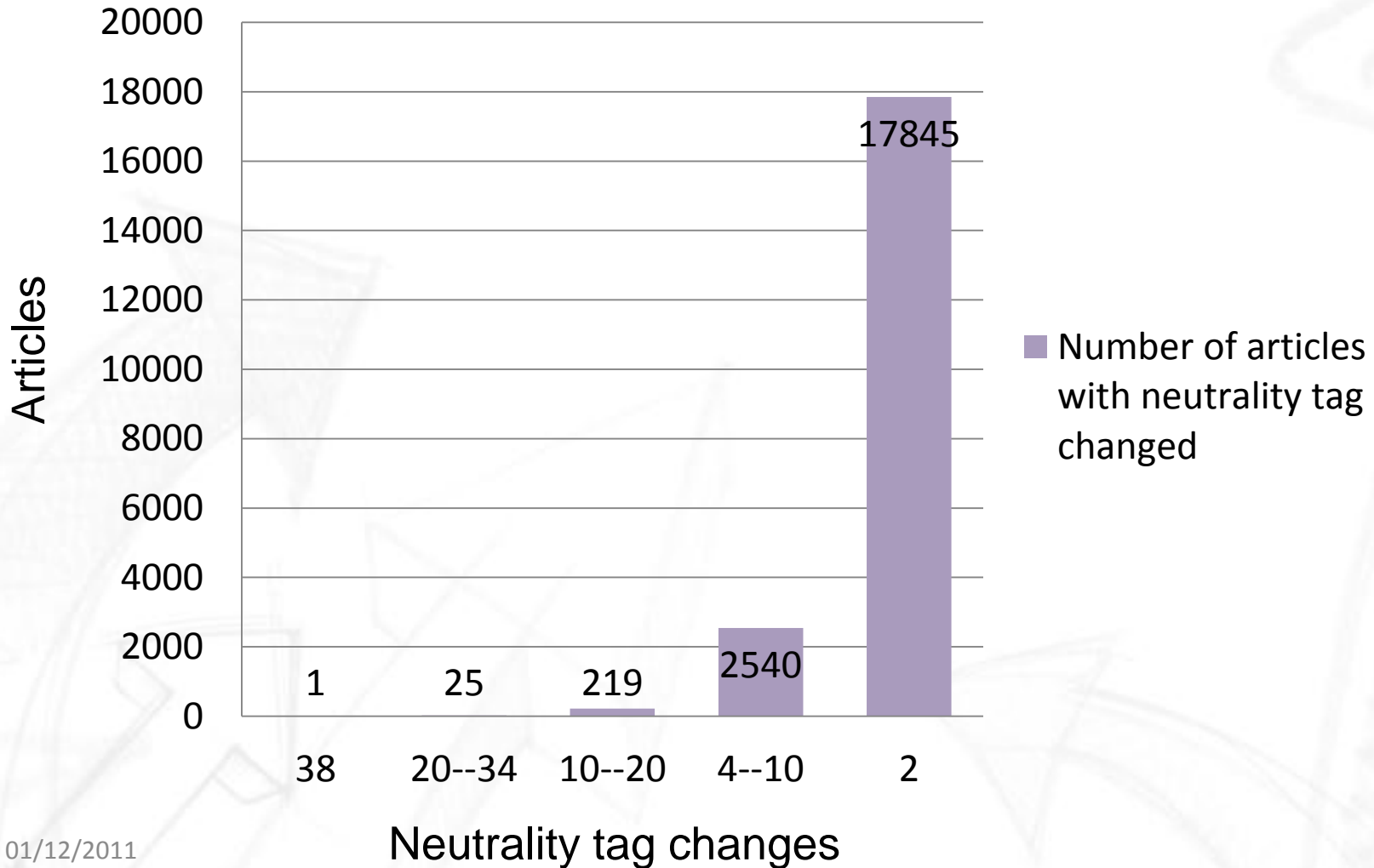
JSI



- Analyzing opinionated Wikipedia articles
 - With the neutrality template set
 - 2 versions of the article (old and new in the dataset)
 - Condition: neutrality template not reset for 7 days after deletion
 - ~ 2GB xml file with opinionated / non-opinionated articles
 - 20,630 opinionated articles
 - 18,719,338 articles in total
- (these counts exclude Wikipedia infrastructure articles)



Wikipedia Neutrality Articles



Top 20 Articles Based on the Neutrality Tag Change

Opinionated Wikipedia articles

September 11 attacks - 38

George W. Bush - 34

Intelligent design - 32

Global warming - 32

Race and intelligence - 28

Circumcision - 28

Armenian Genocide – 28

Macedonians (ethnic group)
- 26

Iraq War - 26

The Holocaust - 24

Muhammad - 24

Jesus - 24

Israel - 24

Arab Israeli conflict - 24

Zionism - 22

Srebrenica massacre - 22

Islamophobia - 22

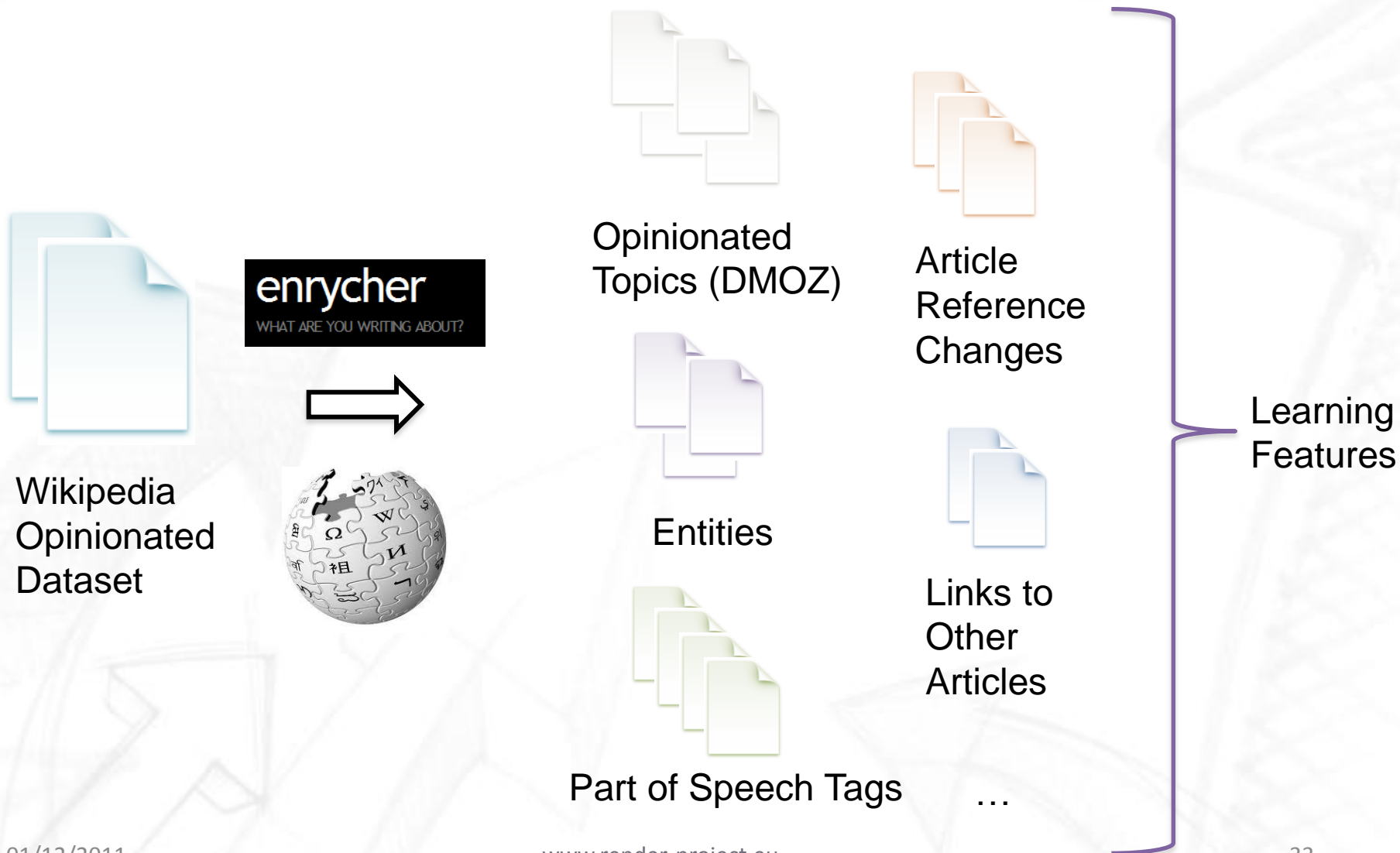
Homeopathy - 22

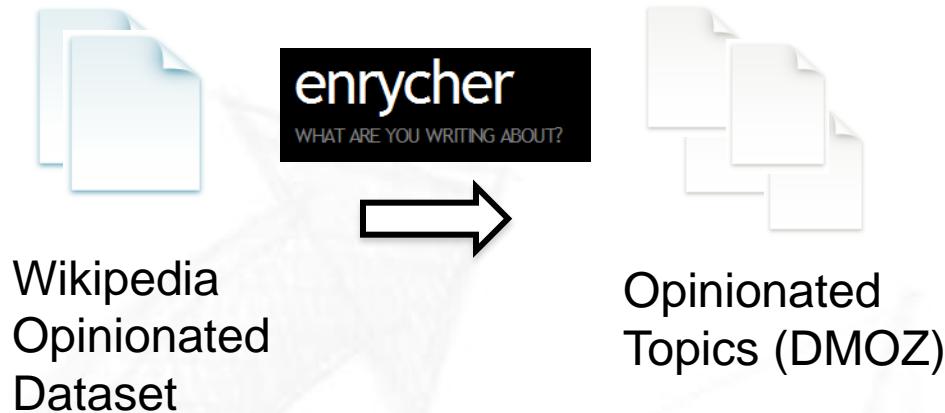
Holocaust denial - 22

Evolution - 22



Learning opinionated Wikipedia articles

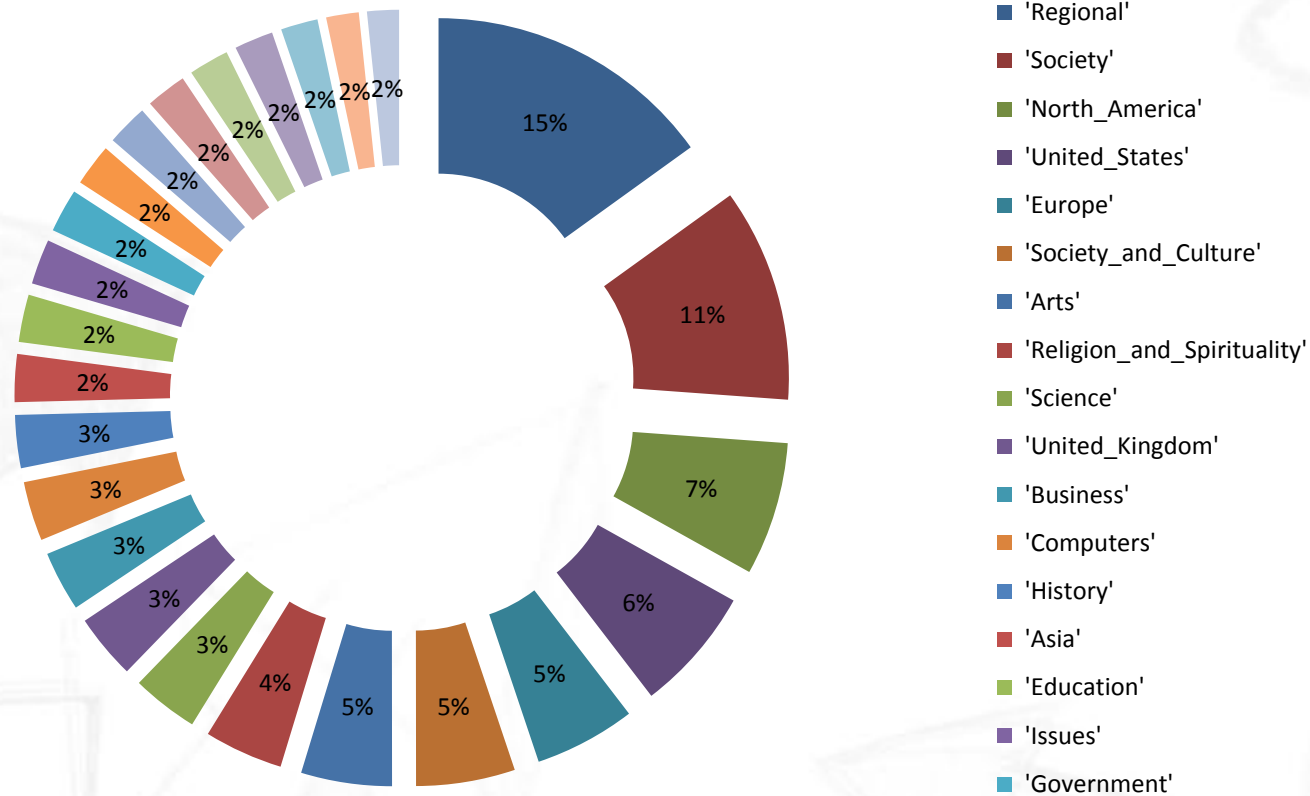




- ~ 10,000 opinionated articles
- Applied Enrycher services:
 - POS tagging
 - DMOZ
 - Topics
 - Keywords



Article Changes for Top Keywords



* Keywords for which the number of article changes is greater than 500



○ **keyword-based changes:**

- article - Megleno-Romanians
- topics - Regional, Europe, Romania, Society_and_Culture, Organizations
- [...] Vlahi is a **disputed** exonym [...]

○ **reference additions:**

- article - Soviet occupation of Romania
- topics - Regional, Europe, Romania
- [...] Sergiu Verona, "Military Occupation and Diplomacy: Soviet Troops in Romania, 1944-1958", Duke University Press [...]



- Next steps in learning opinionated articles:
 - extract the remaining learning features – entities, article references, article links, etc.
- Dealing with scale:
 - process whole Wikipedia articles
 - ~ 30 TB of data
 - extract Wikipedia social network
 - obtain a Wikipedia static and dynamic profile for each contributor/community



Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

Outlook

Detecting bias-inducing editor behavior to generate warnings

KIT



- Relevant use case:
 - UC1: Displaying warnings when detecting patterns of bias
- Intention:
 - Help in understanding and curing behavioral causes for bias
- Needed:
 - Understanding and prediction of socio-technical mechanisms leading to biases





Examples for behavioural-pattern-based warnings presented to users:

- *Concentration*: **98%** of the article were written by **3%** of the active editors in the article. The resulting concentration coefficient is of **9 of 10**. Usual coefficient for similar articles is **5 of 10**. Find out what that means and what you can do to help.
- *Homogeneity*: We detected a **very fractioned** editor structure with **80% of non-vandal edits being reverts** and **3 major editor camps**. Click here for explanation and visualization. Find out what that means and what you can do to help.



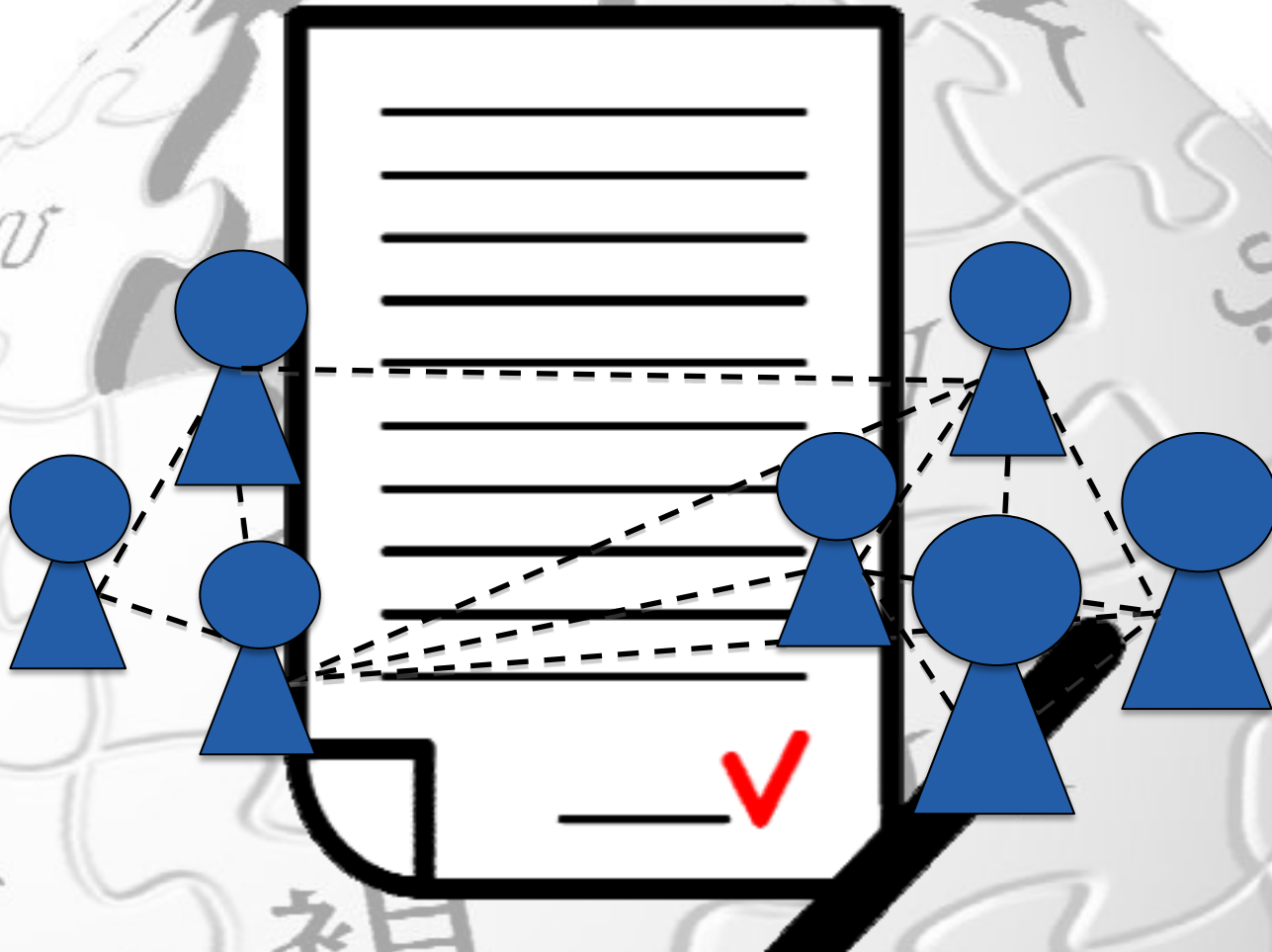
○ Research review

- Identified existing patterns of socio-technical mechanisms potentially influencing bias and diversity
 - Social proof and consolidation
 - Ownership behavior
 - Opinion camps and editor drop-out
 - Lack of boldness and useful conflicts
 - etc.
- see paper “Towards a diversity-minded Wikipedia” and extended literature survey to be published

Identify crucial editing patterns



Example patterns: Opinion Camp drop-out





- Define metrics to find patterns in the data typical for the mechanisms
 - For example: A dense core group of editors in an article's social network structure
- How to get the data for using these metrics?

Reverts as the basis for accurately modeling user behavior in Wikipedia



- Most telling: editors' actions which are *related* to each other
 - Reverting is undoing; contradicting actions perceived as false
 - Inferences possible without knowing meaning
 - Example:

Edit No.	Content	Added/deleted content
1	"zzz"	
2	"zzz yyy"	+ "yyy"
3	"zzz"	- "yyy"

→ Foundation for the search *for* and analysis *of* most of the recurring editing patterns that are typical for biased articles



Simple identity revert method using MD5 hashes

Edit Number	Article content	Words deleted/added (actions taken) by edit	MD5 Hash (simplified)
1	Zero	(ignored for this example)	Hash1

Deficiencies of the state-of-the-art



- Partly reverts exist
- Reverts do not always produce duplicate revisions

Edit Number	Revision content	Words deleted/added (actions taken) by edit	MD5 Hash (simplified)	Detected identic and reverted revisions
1	Zero	(ignored for this example)	Hash1	Like revision 5
2	Zero Apple Banana	+ "Apple" + "Banana"	Hash2	Reverted by revision 5
3	Zero Apple Banana Coconut Date	+ "Coconut" + "Date"	Hash3	Reverted by revision 5
4	Zero Coconut Date	- "Apple" - "Banana"	Hash4	Reverted by revision 5
5	Zero	- "Coconut" - "Date"	Hash1	Like revision 1 → revert of revisions 2,3,4



- A revert is defined by Wikipedia as an action of an editor “undoing the effects of one or more edits” and “(m)ore broadly, reverting may also refer to any action that in whole or in part reverses the actions of other editors.”
- Clear definition, taking into account Wikipedia definition, known intentional behavior & available data:

An edit A is reverted if all of the actions of that edit are completely undone in one subsequent edit B. Edit B has then reverted edit A.

Improved method - implementation



Edit No.	Revision content	Words deleted/added (actions taken) by edit	MD5 Hash (simplified)
1	Zero	(ignored for this example)	Hash1
2	Zero Apple Banana	+“Apple” +“Banana”	Hash2
3	Zero Apple Banana Coconut Date	+“Coconut” +“Date”	Hash3
4	Zero Coconut Date	-“Apple” - “Banana”	Hash4
5	Zero	-“Coconut” - “Date”	Hash1
6	Zero Fig	+“Fig”	Hash5
7	Zero Grape	+“Grape”	Hash6
8	Zero Huckleberry	-“Fig” -“Grape” +“Huckleberry”	Hash7



- Survey evaluation: Accuracy is much higher for new method
 - Significantly less false positives
 - Can accurately distinguish between full and partial reverts
- 12% more reverts detected with the new method than with identity reverts
 - Up to 50% more in short articles
- First revert detection evaluated to work according to the Wikipedia definition and to editors' idea of a revert → better reflects actual behavior and relations → key to precisely modeling the social editing dynamics



Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

Outlook

Mockup demo



Tools to support readers/ editors/ administrators:

- Quality overview of Wikipedia articles for **readers**
- Generation of working lists
 - for Wikipedia **editors** concerning problems of the content
 - Wikipedia **administrators** concerning editor behaviour and interaction

Mock-ups (1) - QAO



WIKIPEDIA The Free Encyclopedia

Article Discussion

Battle of Tippecanoe

From Wikipedia, the free encyclopedia

The **Battle of Tippecanoe** (/tɪpɪˌkɑːnuː/ *TIP-ee-ka-NOO*) was fought on November 7, 1811, between the United States and the Shawnee leader Tecumseh. Tecumseh and his brother Tenskwatawa led a confederacy of Native Americans from various tribes that opposed U.S. expansion. As tensions and violence increased, Governor Harrison marched with an army of about 1,000 men to the confederacy's headquarters at Prophetstown, near the confluence of the Tippecanoe and Wabash rivers. Tecumseh, not yet ready to oppose the United States by force, was away recruiting allies when Harrison's army arrived. Tenskwatawa, a spiritual leader but not a military man, was in charge. Harrison camped near Prophetstown on November 6 and arranged to meet with Tenskwatawa the following day. Early the next morning, warriors from Prophetstown attacked Harrison's army. Although the outnumbered Natives took Harrison's army by surprise, Harrison and his men stood their ground for more than two hours. The Natives were ultimately repulsed when their ammunition ran low. After the battle, the Natives abandoned Prophetstown. Harrison's men burned the town and returned home.

Harrison, having accomplished his goal of destroying Prophetstown, proclaimed that he had won a decisive victory. He acquired the nickname "Tippecanoe", which was popularized in the song "Tippecanoe and Tyler Too" during the election of 1840, when Harrison was elected president. But some of Harrison's contemporaries, as well as some subsequent historians, raised doubts about whether the expedition had been a success. Although the defeat was a setback for Tecumseh's confederacy, the Natives soon rebuilt Prophetstown, and frontier violence actually increased after the battle.

Public opinion in the United States blamed the violence on British interference. This suspicion led to further deterioration of U.S. relations with Great Britain and served as a catalyst to the War of 1812, which began only six months later. By the time the U.S. declared war on Great Britain, Tecumseh's confederacy was ready to launch its war against the United States and embrace an alliance with the British.

Contents [hide]

- Background
- Battle
- Aftermath
- Memorial
- See also
- Notes
- References

Log in / create account

Read View source View history QAO Search

Wikipedia Quality Assessment

General:
 Status: Featured article since 15/05/2010
 Number of edits: 3588
 Number of unique editors: 245
 First edit: 04/03/2006
 Recent edit: 14/10/2011, 17:46
 Number of references: 27 [I more...](#)

Tippecanoe

19th-century depiction by Alonzo Chappel of the final charge that dispersed the Natives^[1]

Date	November 7, 1811
Location	Near modern Battle Ground, Indiana, United States
Result	United States victory ^[2]

Belligerents

Tecumseh's Confederacy	United States
------------------------	---------------

Commanders and leaders

Tenskwatawa	William Henry Harrison
-------------	------------------------

Log in / create account

Wikipedia Quality Assessment

General:

Status: Featured article since 15/05/2010
 Number of edits: 3588
 Number of unique editors: 245
 First edit: 04/03/2006
 Recent edit: 14/10/2011, 17:46
 Number of references: 27

[↑ less...](#)

Further Assessment:

Quality metric I: 0,7
 Quality metric II: 5,3
 Quality metric III: 4,0
 Article feedback score: [View page ratings](#)

[↑ less...](#)

RENDER Analysis:

Fact coverage: [View page ratings](#)
 Neutrality: [View page ratings](#)
 Timeliness: [View page ratings](#)

WIKIPEDIA The Free Encyclopedia

Article Discussion

Battle of Tippecanoe

From Wikipedia, the free encyclopedia

The **Battle of Tippecanoe** (/tɪpɪˌkɑːnuː/ *TIP-ee-ka-NOO*) was fought on November 7, 1811, between the United States and the Shawnee leader Tecumseh. Tecumseh and his brother Tenskwatawa led a confederacy of Native Americans from various tribes that opposed U.S. expansion. As tensions and violence increased, Governor Harrison marched with an army of about 1,000 men to the confederacy's headquarters at Prophetstown, near the confluence of the Tippecanoe and Wabash rivers. Tecumseh, not yet ready to oppose the United States by force, was away recruiting allies when Harrison's army arrived. Tenskwatawa, a spiritual leader but not a military man, was in charge. Harrison camped near Prophetstown on November 6 and arranged to meet with Tenskwatawa the following day. Early the next morning, warriors from Prophetstown attacked Harrison's army. Although the outnumbered Natives took Harrison's army by surprise, Harrison and his men stood their ground for more than two hours. The Natives were ultimately repulsed when their ammunition ran low. After the battle, the Natives abandoned Prophetstown. Harrison's men burned the town and returned home.

Harrison, having accomplished his goal of destroying Prophetstown, proclaimed that he had won a decisive victory. He acquired the nickname "Tippecanoe", which was popularized in the song "Tippecanoe and Tyler Too" during the election of 1840, when Harrison was elected president. But some of Harrison's contemporaries, as well as some subsequent historians, raised doubts about whether the expedition had been a success. Although the defeat was a setback for Tecumseh's confederacy, the Natives soon rebuilt Prophetstown, and frontier violence actually increased after the battle.

Public opinion in the United States blamed the violence on British interference. This suspicion led to further deterioration of U.S. relations with Great Britain and served as a catalyst to the War of 1812, which began only six months later. By the time the U.S. declared war on Great Britain, Tecumseh's confederacy was ready to launch its war against the United States and embrace an alliance with the British.

Contents [hide]

- Background
- Battle
- Aftermath
- Memorial
- See also
- Notes
- References

Log in / create account

Read View source View history QAO Search

Wikipedia Quality Assessment

General:
 Status: Featured article since 15/05/2010
 Number of edits: 3588
 Number of unique editors: 245
 First edit: 04/03/2006
 Recent edit: 14/10/2011, 17:46
 Number of references: 27 [I less...](#)

Further Assessment:
 Quality metric I: 0,7
 Quality metric II: 5,3
 Quality metric III: 4,0
 Article feedback score: [View page ratings](#)

RENDER Analysis:
 Fact coverage: [View page ratings](#)
 Neutrality: [View page ratings](#)
 Timeliness: [View page ratings](#)

Tippecanoe

19th-century depiction by Alonzo Chappel of the final charge that dispersed the Natives^[1]

Date	November 7, 1811
Location	Near modern Battle Ground, Indiana, United States
Result	United States victory ^[2]

Belligerents

Tecumseh's Confederacy	United States
------------------------	---------------

Commanders and leaders

Tenskwatawa	William Henry Harrison
-------------	------------------------

Log in / create account



Biased articles

[This page](#): [Discuss this page](#) - [What does this page mean?](#)

Biased articles for: [Featured articles](#) - [Good articles](#) - [Living people](#)

Biased articles options

Show [50](#) | [100](#) | [250](#) | [500](#) biased articles

Namespace Invert selection Associated namespaces

Category:

Biased articles

[Verotoxin-producing Escherichia coli](#)
[Michael Jackson's health and appearance](#)
[2018 Winter Olympics](#)
...

Neutrality score: X
Neutrality score: Y
Neutrality score: Z
...



Overview

Problem scenarios

Solution including reuse of R&D results

Mockup demo

Outlook

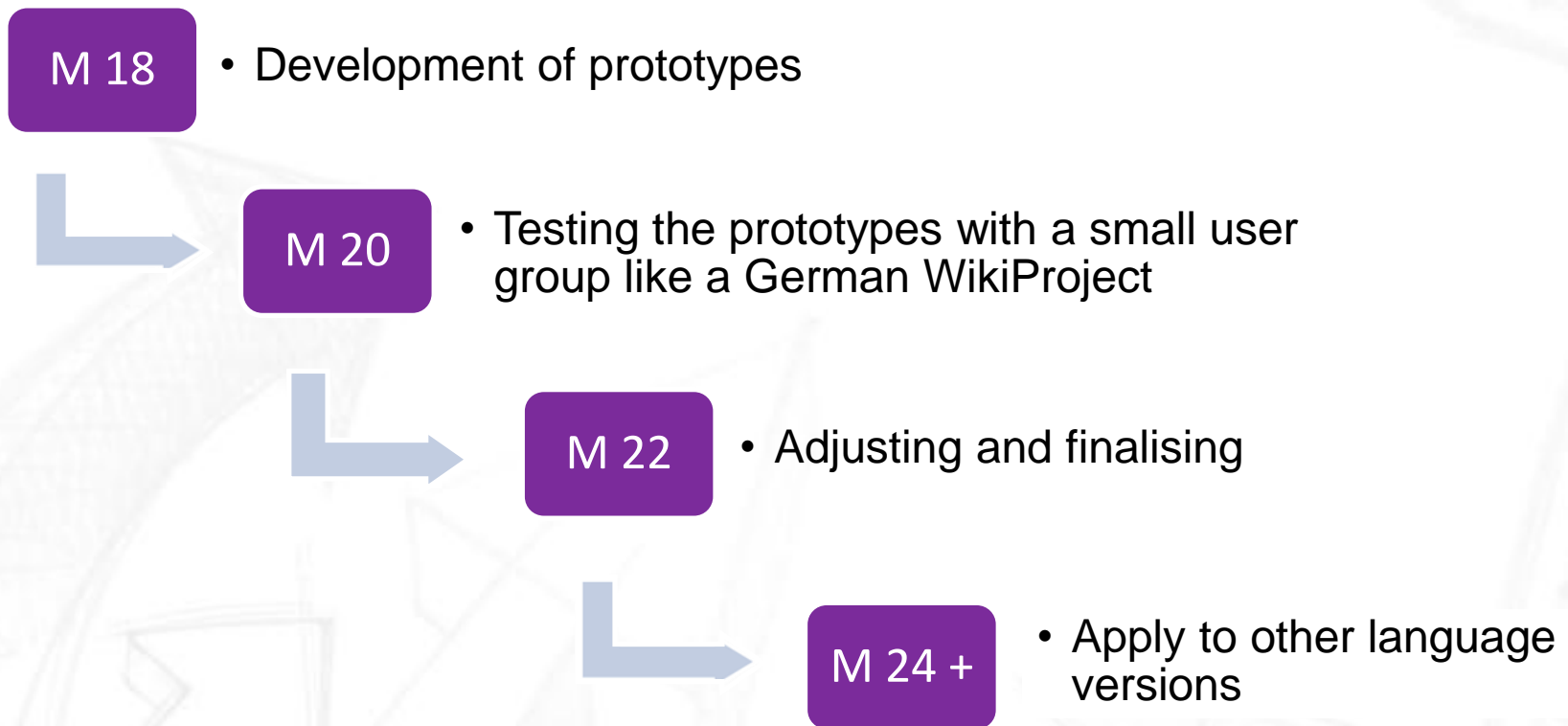
Outlook



- Testing and evaluation of R&D results for Wikipedia
- Development of prototypes for user supporting tools using these results
- Evaluating and testing of these prototypes with Wikipedia users, in the first step of the German community
- Collecting feedback and building up guidelines



Roadmap to develop supporting tools for Wikipedia users:





Questions & comments
Thanks