

Wikitext:

Upcoming changes, Available tools, How you can help

S. Subramanya Sastry, James D. Forrester
Wikimania 2017



WIKIMEDIA
FOUNDATION

[https://etherpad.wikimedia.org/p/
Wikimania2017-wikitext-aug11](https://etherpad.wikimedia.org/p/Wikimania2017-wikitext-aug11)



WIKIMEDIA
FOUNDATION

First things first



WIKIMEDIA
FOUNDATION

Wikitext “parsers”



PHP Parser: default parser (since 2003)

Desktop view, Mobile web, iOS app, Action API

Wikitext “parsers”



Parsoid: alternate parser (since 2012)

VisualEditor, 2017 Wikitext Editor, Flow, Content Translation, Android App,
Linter Extension, Kiwix Offline Reader, Google, REST API

Wikitext “parsers”



Parsing Team @ WMF



WIKIMEDIA
FOUNDATION

Input

Advance wikitext as a language
Easier to write, faster to parse, less error-prone

Wikitext



PHP Parser



HTML

Wikitext



Parsoid



HTML



Make wikitext content easier to analyze

Expose wikitext semantics, no parsing required

Output



“Parsers”

Unify parsers: Same parser for reads as well as edits



Parsing Team Mission

- **Input:** Advance wikitext as a language
 - Easier to write, faster to parse, less error prone
- **Output:** Make wikitext content easier to analyze
 - Expose wikitext semantics in well-specified output
- **Parsers:** Unify parsers
 - Same parser for reads as well as edits

וְשִׁבְעָה יָמִים

Why should you care?

- **Reasoning:** (Structured semantics for wikitext)
 - Improved ability to reason about wikitext (edits, templates)
- **Performance:** (Incremental parsing for high-performance edits)
 - Edits to popular templates won't hurt performance
- **Edits:** (Finer granularity edits → fewer edit conflicts)
 - List items, table cells, quoted text, html tags possibly
 - “Safer” templates, “Easier” template arguments

How this will affect you

- Changes to wikitext (affects editors)
 - Fix edge cases, improve semantics, make slow incremental changes
- Changes to HTML (affects gadget authors, HTML clients)
 - HTML5, semantic markup with a versioned spec
- Changes to parser hooks (affects extension authors)
 - Parser internals not exposed to extensions

Specifics, please!



WIKIMEDIA
FOUNDATION

Examples

- **Finished**
 - Preprocessor changes to fix edge cases in language converter code
- **Ongoing**
 - Replacing Tidy (HTML4) with RemexHTML (HTML5)
 - Identifying broken markup via Linter extension
- **Upcoming**
 - Rendering images with `<figure>` instead of `<div>` **(2017)**
 - New [heredoc](#)-style syntax for multi-template blocks **(2017–2018)**
 - “Balanced” templates **(2018)**

Replacing Tidy



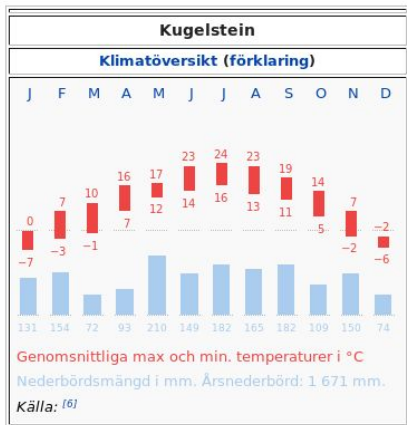
WIKIMEDIA
FOUNDATION

HTML4 → HTML5

- HTML5-based library (RemexHTML) will replace Tidy
- Please read the [FAQ](#) on wiki

Some wikitext might display differently!
Especially broken wikitext / HTML
Or pages that rely on Tidy bugs

och den torraste är mars, med 72 mm nederbörd.^[7]



Kommentarer

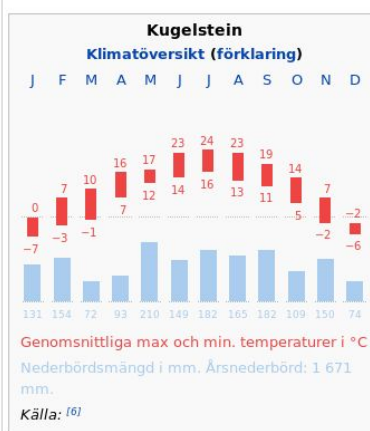
- ↑ Framräknat ur variansen i alla höjduppgifter (DEM 3") från Viewfinder Panoramas, inom 10 km radie.^[2] Mer om algoritmen finns här: [Användare:Lsjbot/Algoritmer](#).
- ↑ Den punkt som svns högst över den lokala horisonten runt platsen, enligt höjduppgifter i GeoNames.^[1]
- ↑ Signifikant i med genomsnittlig förekomst av namngivna sådana på jorden, enligt GeoNames.^[1]

Tidy

Källor

- ↑ [a b c d] [Kugelstein](#) at [GeoNames.Org](#) (cc-by)^[6]; post uppdaterad 2013-07-28; databasdump nerladdad 2015-12-01
- ↑ "Viewfinder Panoramas Digital elevation Model"^[6]. Läst 21 juni 2015.
- ↑ "NASA Earth Observations: Population Density"^[6]. NASA/SEDAC. Läst 30 januari 2016.

och den torraste är mars, med 72 mm nederbörd.^[7]



Kommentarer

- ↑ Framräknat ur variansen i alla höjduppgifter (DEM 3") från Viewfinder Panoramas, inom 10 km radie.^[2] Mer om algoritmen finns här: [Användare:Lsjbot/Algoritmer](#).
- ↑ Den punkt som svns högst över den lokala horisonten runt platsen, enligt höjduppgifter i GeoNames.^[1]
- ↑ Signifikant fler inom 20 km radie jämfört med genomsnittlig förekomst av namngivna sådana på jorden, enligt GeoNames.^[1]

Källor

- ↑ [a b c d] [Kugelstein](#) at [GeoNames.Org](#) (cc-by)^[6]; post uppdaterad 2013-07-28; databasdump nerladdad 2015-12-01
- ↑ "Viewfinder Panoramas Digital elevation Model"^[6]. Läst 21 juni 2015.
- ↑ "NASA Earth Observations: Population Density"^[6]. NASA/SEDAC. Läst 30 januari 2016.
- ↑ "NASA Earth Observations: Land Cover Classification"^[6]. NASA/MODIS. Läst 30 januari 2016.
- ↑ Peel, M C; Finlayson, B L; McMahon, T A (2007). "Updated world map of the Köppen-Geiger climate classification"^[6]". *Hydrology and Earth System Sciences* 11 (6): 403–414. doi:10.5194/hess-11-403-2007.

RemexHTML

↑ [Set Index](#)^[6]. NASA. Läst 30

↑ [1 month - TRMM](#)^[6].

NASA/Tropical Rainfall Monitoring Mission. Läst 30 januari 2016.

Snapshoted via [?action=parsemigration-edit](#)



WIKIMEDIA
FOUNDATION

V • R • U	Departmani Francuske [sakrij]
	<ul style="list-style-type: none"> 01 Ain • 02 Aisne • 03 Allier • 04 Alpes-de-Haute-Provence • 05 Hautes-Alpes • 06 Alpes-Maritimes • 07 Ardèche • 08 Ardennes • 09 Ariège • 10 Aube • 11 Aude • 12 Aveyron • 13 Bouches-du-Rhône • 14 Calvados • 15 Cantal • 16 Charente • 17 Charente-Maritime • 18 Cher • 19 Corrèze • 2A Corse-du-Sud • 2B Haute-Corse • 21 Côte-d'Or • 22 Côtes-d'Armor • 23 Creuse • 24 Dordogne • 25 Doubs • 26 Dôme • 27 Eure-et-Loir • 28 Gard • 29 Gironde • 30 Gers • 31 Haute-Garonne • 32 Gers • 33 Hérault • 34 Ille-et-Vilaine • 36 Indre • 37 Indre-et-Loire • 38 Isère • 39 Jura • 40 Landes • 41 Loir-et-Cher • 42 Loire • 43 Haute-Loire • 44 Loire-Atlantique •

Tidy

V • R • U	
Metropolitanska Francuska	01 Ain • 02 Aisne • 03 Allier • 04 Alpes-de-Haute-Provence • 05 Hautes-Alpes • 06 Alpes-Maritimes • 07 Ardèche • 08 Ardennes
Prekomorski departmani	971 Guadeloupe • 972 Martinique • 973 Francuska Gvajana • 974 Réunion • 976 Mayotte

RemexHTML

Snapshotted via [?action=parsemigration-edit](#)



How you can help

- Fix High Priority errors exposed by [Extension:Linter](#)
 - Exposed via `Special:LintErrors` and the [REST API](#)
 - `Extension:ParserMigration` will help you verify fixes
 - `?action=parsermigration-edit` instead of `?action=edit`
- Status tracking
 - [Dashboard](#) (<https://tools.wmflabs.org/wikitext-deprecation/>)
 - Weekly stats published at [mw:Parsing/Replacing_Tidy/Linter/Stats](#)
 - Weekly pixel-diff run results at <http://mw-expt-tests.wmflabs.org/>

TL; DR



WIKIMEDIA
FOUNDATION

Summary: Get Involved

- Incremental changes to wikitext, HTML, parser APIs
 - Pages, templates, bots, gadgets, extensions might need fixing
- Some changes:
 - Tidy (HTML4) → RemexHTML (HTML5) ([FAQ](#))
 - `<figure>` markup for images ([T118517](#))
 - Improved syntax for multi-template content blocks ([T114432](#))
 - Balanced templates ([T114445](#))
- **Tools:** Linter, ParserMigration, Pixel diffs, Dump greps

https://www.mediawiki.org/wiki/Parsing/Get_involved

DISCUSS!
THANK YOU!



WIKIMEDIA
FOUNDATION

Backup Slides



WIKIMEDIA
FOUNDATION

Here doc syntax (not final)

Now: Four “unconnected”
transclusions

```
{{tablestart}}
```

```
{{tablerow|1}}
```

```
{{tablerow|2}}
```

```
{{tableend}}
```

One Proposal: Single
“connected” transclusion

```
{{<<MyTable}}
```

```
{{tablerow|1}}
```

```
{{tablerow|2}}
```

```
{{>>MyTable}}
```

Benefits

- Eliminates escaping hazards by hoisting content out of template argument position
- Supports balanced templates goal

Balanced Templates

- Add (type) annotation to templates indicating use context

```
{{#balance}}, {{#balance:inline}}, {{#balance:block}},  
{{#balance:attribute}}, {{#balance:string}} ....
```
- Type enforced by parsing to:
 - Well-formed & balanced DOM → unclosed tags closed
 - Attribute → Appropriate escaping without explicit nowikis
 - String → Appropriate escaping without explicit nowikis
 -

Benefits

- Predictable rendering across all uses
- Markup errors are constrained to limited context (table, list, paragraph, etc.)
- Enables higher performance via incremental parsing