#### Dictionary Reversal

Building a German-Japanese Dictionary from Japanese-German Data, Using Vocabulary Analysis, User Cooperation and Approaches from Natural Language Processing

Ulrich Apel

Information Systems Architecture Science Research Division National Institute of Informatics, Tokyo

#### The Presentation

- \* Personal background
- \* The WaDoku Project
  - example applications of the data
  - data structure and user collaboration
- \* Dictionary reversal
  - preliminary project and problems
  - data preparation, additional data and expected results
- \* Conclusion

# Personal Background

- \* Study of japanology (Japanese studies), sociology and ethnology in Munich, Germany
- \* Doctoral thesis on Japanese futures studies (futurology), research at Ōsaka university
- \* Lack of up-to-date Japanese-German dictionary lead to WaDoku Dictionary Project
  - "newest" dictionary from 1980
  - comprehensive dictionary from 1952 (original version from 1937 – old kanji and kana·zukai)
  - situation of Japanese-English hardly better

# WaDoku Project

- \* Very comprehensive Japanese-German electronic dictionary
- \* FileMaker database with more than 100,000 headwords, about 250,000 records
- \* With assumed 40 headwords per page, size corresponds to 2,500 printed pages
- \* Highly flexible and easily accessible data
- \* Look-up of German is possible (Japanese entry with pronunciation)

#### User Access to the Dictionary

- \* 1999
  - FileMaker runtime as download
  - online version on server at Ōsaka
- \* Web page <WaDoku.de>, better user collaboration
  - 2,000–4,000 visits/day
  - o 25,000-45,000 page views/day
- \* Data used in JMDict, Papillon Project, Reading Tutor, PopJisyo, W3dict etc.
- \* Other formats: edict, EPWING, Firefox sidebar etc.

# EPWING format in Jamming



# LingoPad on Windows

the water and the second the states in the Low day to the second with the second a - wine a

DE JP	🖮 🤹 • 🏋 🗖   🗐 🛛	victionaries 🔻 📝 Program 🔻 🕑 Info 🔻
鳥居	▼ ≓ Japanese	German
<u>鳥居</u> 鳥居前町 鳥屋につく 鳥山明 鳥島 鳥打ち 鳥打ち帽 鳥野帽 鳥撃 鳥撃 鳥撃 鳥撃 り 島 二 の し の し の し の し の し の し の し の し の し の	<ul> <li>● 鳥居</li> <li>●</li> </ul>	Rel., Archit. Torii n (Tor vor Shintō-Heiligtümern; wird gebildet aus zwei durch zwei Querbalken verbundene senkrechte Holzpfeilern; früher meist aus Holz; jetzt mitunter auch Stein, Metall oder Beton)
鳥獣が餌をあさる	▼	鳥居 [とりい] Pronunciation

# Moji Sidebar for Firefox

is an energy the transmitter started a . Advant

The latter destine in the washing man and the second stand their

0	Moji	_	1 - 57 1 3	本文ノート 編集 漫歴
言 Wort 漢	Kanji Mehr		and when	
Search Word	• • Meaning	• 🖽	A PHE TO	鳥居
鳥居	Torii		Jes C	
百科事典	Enzyklopädie	-	ウィキペディア	出典: フリー百科事典『ウィキペディア
神域	Schreinbezirk		フリー百科事典	
区画	Bezirk	X	ナビゲーション	<b>鳥居(とりい</b> )は、神社などにおいて、
λD	Eingang	Ŧ	= メインページ	「門」である。御陵や寺院に建てられ
自日	0		■ コミュニティ・ポータル	鳥居を図案化したものになっている。
局店 [とり	[4]		■ 最近の出来事	
{Rel., Archit.}	; Torii; (wörtl. etwa		■ 最近更新したページ	<b>目次</b> [非表示]
"Vogelsitz"; To	r vor Shintö-Heiligtümen ei senkrechten Holznfeile	1;	■ おまかせ表示	1 巻 5 古
die durch zwei	Querbalken verbunden si	nd)	= アップロード (ウィキメ	
	-	-	ディア・コモンズ)	2 起源
			<ul> <li>ウィキペディアに関する</li> </ul>	3 形式
			お問い合わせ	3.1 神明鳥居
			ヘルプ	32 明神鳥居

= ヘルプ

4 種類(材料)

# POPjisyo/POP辞書



# wadokujt.w3dict.com

the storest of fairing in the Low side The States with the stranged a - Barrage

日本語 English Deutsch \*

# wadokujt. w3dict.com

Lesehilfe Japanisch-Deutsch

examples

about

Mala stor waster

contact

#### examples

dictionary-in-text examples. Move mouse pointer over the text.

A	Article 9, Constit	ution of Japan
1	.日本国民は、正義	と秩序を基調とする
Ξ	国際平和 を 誠実に 希求	し、国権の発動た
2	国際平和	威嚇又は武力の行
5	こくさいへいわ	する手段としては、
ź	internationaler	
	Frieden {m}.	
2	kokusai∙ <heiwa></heiwa>	ため、陸海空軍そ
σ		保持しない。国の
3	を戦権は、 これを認め	oない。

# wadoku.de

in the second to the second with the stand of a fairing

BEALTH AND A TOT DO TO MAN STATISTICS TO THE STATE STATE



2 items found	, displaying	all items.
---------------	--------------	------------

Nr:	Japanisch 🗢	Lesung 🗢	Deutsch ¢
1	鳥居	とりい	{Rel., Archit.} Torii {n} (Tor vor Shintō-Heiligtümern; wird gebildet aus zwei durch zwei Querbalken verbundene senkrechte Holzpfeilern; früher meist aus Holz; jetzt mitunter auch Stein, Metall oder Beton).
2	鳥居前町	とりいまえまち	um bzw. vor einem Schrein entstandene Stadt {f} (wie z.B. lse).

# User collaboration

- \* Comments on entries on the web page
- \* Proposals for new entries and corrections
- \* Corrections through editors
- \* Forum
- \* Wiki
- \* Maintenance of web page by volunteers
- \* NPO WaDoku e.V.
- \* Good basis for organic growth of the project

# Data and Possible Representation

DaTimeS	DaJapanisch	DaLesung	DaDeutsch	DaWortart	DaEintragsebene	DaUE_Art	DaUntereintrMidashigo	DaRomajiBe
11.04.2007 19:15:18	飲酒 [1]	いんしゅ [1]	( <pos: n.="">) Trinken<gen.: n=""></gen.:></pos:>	名	HE			いんしゅ
	飲酒する	いんしゅする	trinken ( <expl.: alkohol="">).</expl.:>	サ変自	飲酒[1]	Abl. mit <umschr :<="" td=""><td>&lt;飲酒&gt;する</td><td>&lt;いんしゅ&gt;</td></umschr>	<飲酒>する	<いんしゅ>
12.06.2007 20:13:32	飲酒運転者	いんしゅうんてんしゃ	unter Alkoholeinfluss stehender	名	飲酒[1]	Komp. Anf.	<飲酒>運転者	<いんしゅ>・う
	飲酒狂	いんしゅきょう	[1] Alkoholismus <gen.: m="">.</gen.:>	名	飲酒[1]	Komp. Anf.	<飲酒>狂	<いんしゅ>・き
	飲酒恐怖症	いんしゅきょうふしょ -	{ <dom.: med.="">} Dipsophobie<gen.:< td=""><td>名</td><td>飲酒[1]</td><td>Komp. Anf.</td><td>&lt;飲酒&gt;恐怖症</td><td>&lt;いんしゅ&gt;・き</td></gen.:<></dom.:>	名	飲酒[1]	Komp. Anf.	<飲酒>恐怖症	<いんしゅ>・き
	飲酒検査	いんしゅけんさ	Alkoholtest <gen.: m="">.</gen.:>	名	飲酒[1]	Komp. Anf.	<飲酒>検査	<いんしゅ>・k
	飲酒検知器	いんしゅけんちき	Alkoholdetektor <gen.: m="">.</gen.:>	名	飲酒[1]	Komp. Anf.	<飲酒>検知器	<いんしゅ>・k
	飲酒癖	いんしゅへき	Trunksucht <gen.: f="">; Hang<gen.:< td=""><td>名</td><td>飲酒[1]</td><td>Komp. Anf.</td><td>&lt;飲酒&gt;癖</td><td>&lt;いんしゅ&gt;・^</td></gen.:<></gen.:>	名	飲酒[1]	Komp. Anf.	<飲酒>癖	<いんしゅ>・^
	飲酒量	いんしゅりょう	konsumierte Alkoholmenge <gen.: f="">.</gen.:>	名	飲酒[1]	Komp. Anf.	<飲酒>量	<いんしゅ>・り
	飲酒にふける	いんしゅにふける	sich dem Trinken hingeben; trinken.	下一他	飲酒[1]	Verwendung sheispiel	<飲酒>にふける	<いんしゅ>
	飲酒を慎む	いんしゅをつつしむ	mäßig trinken; beim Trinken Maß	五他	飲酒[1]	Verwendung sheispiel	<飲酒>を慎む	<いんしゅ>

- inshu 飲酒 (N.) Trinken n (von Alkohol).
- →~する trinken (Alkohol).
- △ ~癖 ~*heki* Trunksucht *f*, Hang *m* zum Alkoholismus. II ~家 ~*ka* (N.) Trinker *m*; Zecher *m*; Säufer *m*. II ~検知器 ~*kenchiki* Alkoholdetektor *m*. II ~検査 ~*kensa* Alkoholtest *m*. II ~狂 ~*kyō* ① Alkoholismus *m*. ② Alkoholiker *m*. II ~恐怖症 ~*kyōfushō* MED. Dipsophobie *f*; Angst vor dem Trinken. II ~量 ~*ryō* konsumierte Alkoholmenge *f*. II ~運転 ~*unten* (N.) Trunkenheit *f* am Steuer; Fahren *n* unter Alkoholeinwirkung. II ~運転者 ~*unten·sha* unter Alkoholeinfluss stehender Fahrer *m*; betrunkener Fahrer *m*.
- ☆~にふける sich dem Trinken hingeben; trinken. II ~を慎む ~ o tsutsushimu mäßig trinken; beim Trinken Maß halten.

# Highly Flexible Data

- Mostly automated Romanization in several flavours (segmentation, mark-up for Japanese pitch accent nasalisation of velar plosive and devocalization, capitalisation, joshi は wa and へ he)
- \* Parts of speech, conjugation, gender of nouns
- \* Domains, definitions, comments on usage, etymology, references to synonyms and antonyms
- \* Main entries and corresponding compounds, derivations, examples of usage & example sentences

# Dictionary for MekiMekiDoitsugo

- \* Project for "Germany in Japan 2005/06"
- \* German course for mobile phones
- \* Automated German pronunciation in katakana
- \* Simple reversible dictionary from Japanese-German word pairs for 10,000 most common German words
- \* Added data on pronunciation, part of speech etc.
- \* Short comings of this approach became obvious (uncommon Japanese writings, obsolete entries, missing information for Japanese users etc.)

# Problems for Dictionary Reversal

- \* Grammatical distance no word-to-word correspondence (noun translated as adjective)
- \* Cultural distance no semantical one-to-one correspondence (definition instead of translation)
- *\** Different needs of German and Japanese users
   *o* information about cultural background
  - no need for Romaji transcription, or gender of German entries and inflection
  - Research necessary, e.g. questionnaires
- \* Complex structure of WaDoku data

# Problems of Data and Structure

- \* Many orthographical variants Japanese entries
- \* Appropriate translation even on expense of the quality of the German (e.g. many nominalisation or not typical German)
- \* Explanations on cultural background are difficult to process
- \* Important information for Japanese users is missing — German pronunciation, gender, no information on inflection or valency of verbs etc.

# Preparation of WaDoku Data

- \* Mark-up for basic meanings of Japanese entry
- \* Mark-up for main meanings of Japanese entry (likely to be easily reversible)
- \* Mark-up for reversible translations "harmonic dictionary (Kumiko Tanaka-Ishii)
- \* Mark-up for non-reversible translations

# Building a Raw Reversed Dictionary

- \* Create records of one "official" Japanese writing, one reading and one German translation
- \* Delete cultural and grammatical explanation on Japanese entry (POS, definitions, explanations etc.)
- \* Replace actual German word forms with dictionary form if necessary (e.g. singular instead of plural)
- \* Add available information about German entry (pronunciation, inflection, gender etc.)

# Improvement of the Raw Dictionary

#### \* User cooperation

- give input about their needs for the dictionary
- rank translation candidates according to suitability and group sub-meanings
- add and correct translations
- \* Analyses of the entries and NLP
  - domains with good or bad one-to-one correspondence (modern sciences vs. religion)
  - mark up unlikely candidates (e.g. archaic usage)
  - probabilities through e.g. frequencies in corpuses

# **Expected Results**

- \* Raw dictionary
  - usable dictionary in relatively short time
  - decent size but relatively mediocre quality
- \* Improvements through user cooperation
  - o good experience with WaDoku.de
  - several iterations necessary
- \* Improvements through analyses and NLP
  - in parts very promising
  - in parts uncertain

# Conclusion

- \* WaDoku dictionary one of the most important electronic tools for German research on Japan
- \* DokuWa dictionary might play a similar role for Japanese research on Germany
- \* Developed method could be applied to other dictionary projects
- \* Data probably basis for "harmonic dictionary"
- Relatively big benefit through relatively small effort

#### Thank You for Your Attention!

a tor wast's me we corperations in on the

and the transmithe states of a - Adverte