

# Open Source-MÜ-Systeme

## Seminar «Maschinelle Übersetzung» (HS 2011)

Simon Hafner & Hernani Marques

30. November 2011

# Referatsübersicht

- 1 Problemstellung
  - Übersetzen von und nach Minderheitensprachen
  - Closed Source vs. Open Source Software
- 2 Herausforderungen & Lösungsansätze
- 3 Zwei OSS-MÜ-Systeme: Apertium und Moses
  - Apertium
  - Moses
- 4 Fazit
- 5 Fragen

# Minderheitensprachen

- Sprachen einer (kleinen) (Sub-)Kultur oder einer kleinen (geografischen) Region
- Beispiele: Klingonisch, Quechua, Baskisch, Esperanto

# Problem: Bei vielen MÜ-Systemen keine Unterstützung

- Grosse/bekannte/kommerzielle MÜ-Systeme orientieren sich an Nachfrage/Rendite
- Datenmaterial für Minderheitensprachen ist nur spärlich vorhanden

# Open Source Software (generell)



Abbildung: Open Source Software (Logos)<sup>1</sup>

**Lizenzen** GPL, LGPL, CC, MIT/BSD/Apache, Shared Source u. a.  
(libertär/viral, liberal, zweckbeschränkt)

**Software** Betriebssysteme, Browser, Spiele u. a.; auch:  
MÜ-Systeme :)

---

<sup>1</sup>URL: <http://cdn.lostintechology.com/wp-content/uploads/2010/12/opensource.jpg> (30.11.2011)

# Closed Source Software (generell)

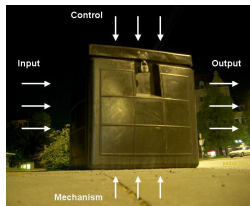


Abbildung: Closed Source Software (als Blackbox verstanden)<sup>2</sup>

*Lizenzen* EULA, Mehrplatzlizenzen, Shareware, Freeware u. a.

*Software* Betriebssysteme, Browser, Spiele u. a.; auch:  
MÜ-Systeme :)

---

<sup>2</sup>URL: [http://juid.de/images/Kaleidoskop/The\\_black\\_box.jpg](http://juid.de/images/Kaleidoskop/The_black_box.jpg) (30.11.2011)

# Open Source Software: Vorteile (1)

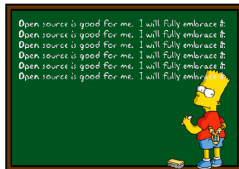


Abbildung: Bart: \*embraceOpenSource\*<sup>3</sup>

- *Programmcodes* / *Daten* sind kostenfrei (üblich)
- *Änderung/Erweiterung* des Programmverhaltens und der -funktionalität möglich
- *Nachvollzug* des Outputs möglich

---

<sup>3</sup>URL: <http://toomanytabs.com/blog/wp-content/uploads/2010/11/Bartopen.gif> (30.11.2011)

# Open Source Software: Vorteile (2)

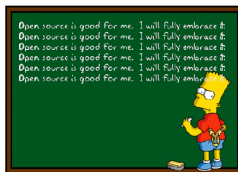


Abbildung: Bart: \*embraceOpenSource\*<sup>4</sup>

- *Kultur* des Teilens / kollektive Entwicklung
- *Partizipation* ist möglich und erwünscht

---

<sup>4</sup>URL: <http://toomanytabs.com/blog/wp-content/uploads/2010/11/Bartopen.gif> (30.11.2011)



# Herausforderungen

- Minderheitensprache implementieren
- Kein kommerzieller Support
- Sparse-Data Problem
- stark agglutinierende Sprachen

# Herausforderung (1)

**Problem** Kein kommerzieller Support

**Lösung** Auf den Schultern von Giganten

**Lösung** OSS Systeme

## Herausforderung (2)

**Problem** Sparse-Data Problem

**Problem** stark agglutinierende Sprachen

**Lösung** Regeln implementieren

**Lösung** Zirkuläre Implikation [Forcada2006]

# Apertium: Merkmale (1)



Abbildung: Apertium-Logo<sup>5</sup>

- Entwicklungsschwerpunkt in Spanien
- GPL-lizenziert

---

<sup>5</sup>URL: <http://www.eamt.org/corporate/logos/apertium.png> (30.11.2011)

# Apertium: Merkmale (2)



Abbildung: Apertium-Logo<sup>6</sup>

- RBMT
- Unterstützung für 28 Sprachpaare

---

<sup>6</sup>URL: <http://www.eamt.org/corporate/logos/apertium.png> (30.11.2011)

# Technologien (Shallow-Transfer)

- Für syntaxähnliche Sprachen
- 1-Pass-Modell

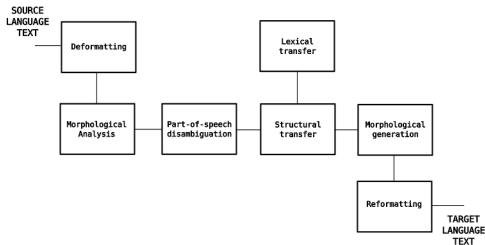


Abbildung: Apertium-Module<sup>7</sup>

---

<sup>7</sup>URL:

[http://wiki.apertium.eu/images/thumb/9/99/Figure\\_1\\_dibuix.svg/768px-Figure\\_1\\_dibuix.svg.png](http://wiki.apertium.eu/images/thumb/9/99/Figure_1_dibuix.svg/768px-Figure_1_dibuix.svg.png)  
(30.11.2011)

# Technologien (Advanced Transfer)

- Für syntax~~un~~ähnliche Sprachen
- 3-Pass-Modell: *Chunking* - *Interchunking* - *Postchunking*

# Technologien (XML)

- Linguistische Daten werden in der XML gehalten
- Interoperabilität ist damit gesichert

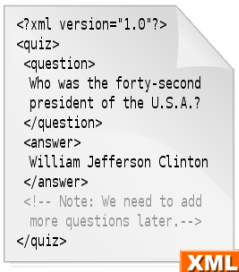


Abbildung: XML-Beispiel<sup>8</sup>

<sup>8</sup>URL: <http://upload.wikimedia.org/wikipedia/commons/thumb/6/68/XML.svg/275px-XML.svg.png>  
(30.11.2011)



# Phrase-Based Alignment



Abbildung: Phrases<sup>9</sup>

<sup>9</sup>URL: <http://www.statmt.org/moses/img/waiph5.png> (30.11.2011)

# Phrase-Based Applied

Maria	no	daba	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap	by			green witch	
	no		slap		to the			
	did not give				to			
					the			
			slap			the witch		

Abbildung: Translations<sup>10</sup>

<sup>10</sup>URL: <http://www.statmt.org/moses/img/translation-options.png> (30.11.2011)

# Beam Search

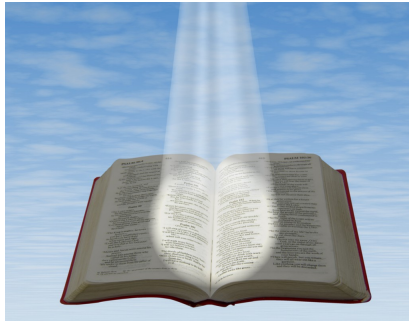


Abbildung: Beam Search<sup>11</sup>

---

<sup>11</sup>URL: <http://seo.mindandmouth.com/wp-content/uploads/2010/03/bible-illuminated-by-beam-of-light1.jpg?f8021c>

(30.11.2011)

# Beam Example

Maria	no	daba	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not		a slap	by			green witch	
	no		slap		to the			
	did not give				to			
					the			
			slap			the witch		

Abbildung: Translations<sup>12</sup>

<sup>12</sup>URL: <http://www.statmt.org/moses/img/translation-options.png> (30.11.2011)

# Beam Applied

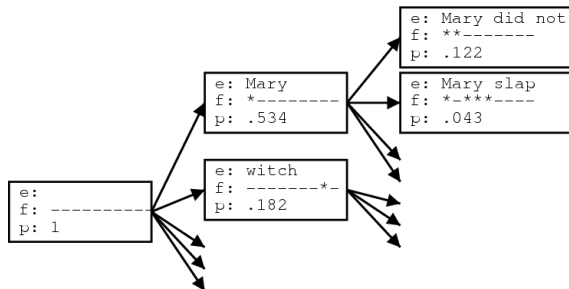


Abbildung: Translations<sup>13</sup>

<sup>13</sup>URL: <http://www.statmt.org/moses/img/beam.png> (30.11.2011)

# Beam Search Data

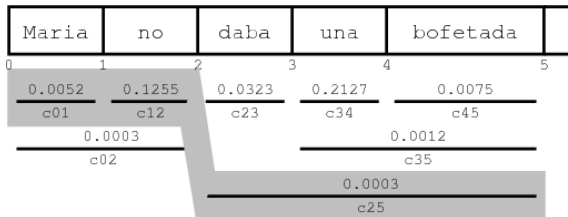


Abbildung: Future Cost<sup>14</sup>

<sup>14</sup>URL: <http://www.statmt.org/moses/img/phrase-beam-future-cost.png> (30.11.2011)

# Factored

- Multi-Dimensionale Daten
- z.B. Morphologie

# Factored Example

*jllegh'egh*

*[SG][1PERS][NONE][VERB]legh[oneself]*



# Lektüreempfehlungen

[Norvig1992] Norvig, Peter: *Paradigms of artificial intelligence programming: case studies in common LISP*. S.195-196. San Fransisco: Morgan.

[Forcada2010] Forcada Mikel L. et al.: *Documentation of the Open-Source Shallow-Transfer Machine Translation Platform Apertium*.

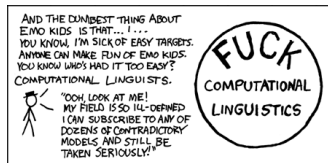
URL: <http://xixona.dlsi.ua.es/~fran/apertium2-documentation.pdf> (30.11.2011)

[Forcada2006] Forcada, Mikel L.: *Open-source machine translation: an opportunity for minor languages*. In: Strategies for developing machine translation for minority languages (5th SALT MIL workshop on Minority Languages).

URL: <http://www.dlsi.ua.es/~mlf/docum/forcada06p2.pdf> (30.11.2011)

# Auswirkungen: Einsatz von OSS-MÜ-Systemen

- Minderheitensprachen können durch Open Source-MÜ-Systeme mehr Verbreitung finden
  - RBMT: Apertium
  - SBMT: Moses
- Sprechende von Minderheiten können sich stärker ihrer eigenen Sprache bedienen

Dalegh'a<sup>16</sup>Abbildung: XKCD - Nr. 114<sup>15</sup>

---

<sup>15</sup>URL: <http://xkcd.com/114/> (30.11.2011)

<sup>16</sup>engl. "Do you see it?". URL: <http://en.wikibooks.org/wiki/Klingon/Grammar/Questions> (30.11.2011)