# Tools for Translation State of the Practice, Established Tools and Capability Gaps

**Jennifer DeCamp**

MITRE Corporation

jdecamp@mitre.org

# Agenda

- Fundamental Questions
- Types of Translation
- Translation Workflows and Current Tools
- Gaps
- Tolerance/Preferences
  - Human Translator
  - Post-Editor
  - MT User
- Conclusions
- Recommendations for MT Summit

# Fundamental Questions

1. How can you get better accuracy in MT/HT?

   Communicative value (Tapling 2008)

   Fluency, fidelity, etc.

2. How can you make a C-level translator into a B-level translator?

3. How can you make human translators more productive?

# Types of Translation

- Translation
- High Quality Translation
  - Literature, marketing material, important briefings
  - Nuance, style
- Gists/summaries
- Getting information to answer specific questions
- Sorting—figuring out the language, subject, and difficulty in order to route the material
- Multimedia
  - Television, etc.
- Embedded MT
- Multi-language source material
- Different dialects and registers
- MT embedded in chat and search tools
- Term translation
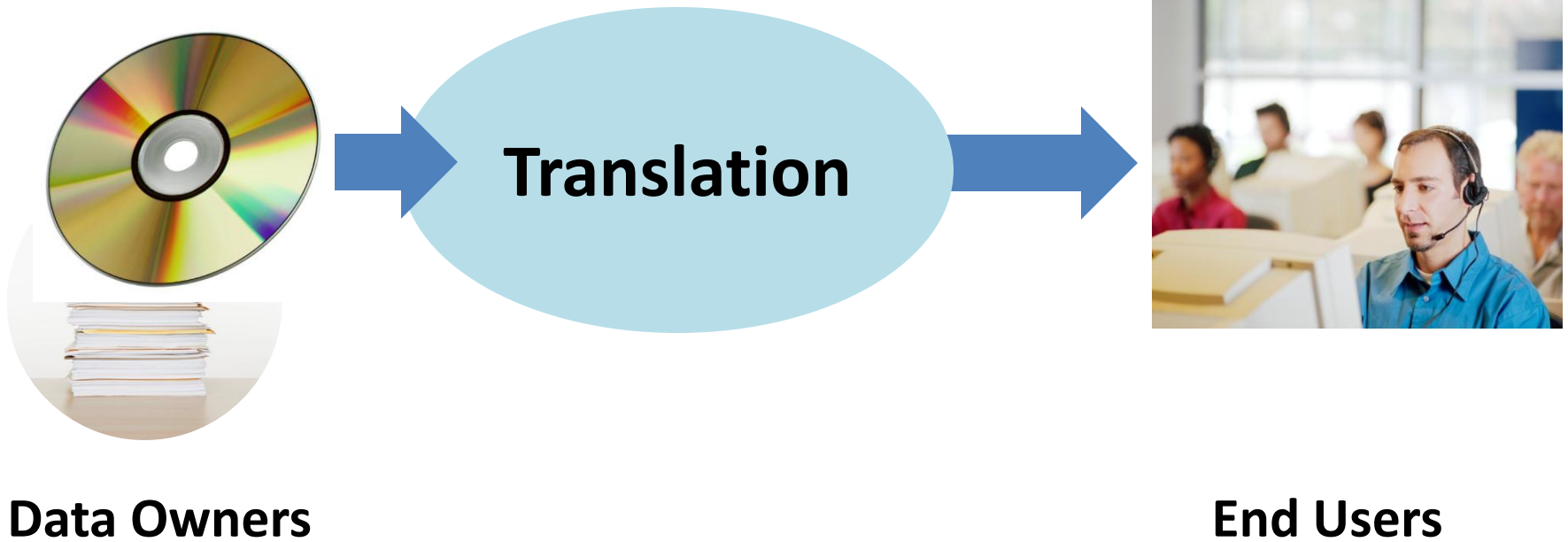- Speech-to-speech translation

# Some Major Distinctions

- Publisher centric vs. User centric
- Human Translation (HT) vs. Machine Translation (MT)
  - Increasingly mixed
- Cautious vs. whatever

# Combinations

- Publisher-Centric
  - HT
  - MT
  - Combination
- User-Centric
  - HT
  - MT
  - Combination
  - With user tools?

# Publisher-Centric HT



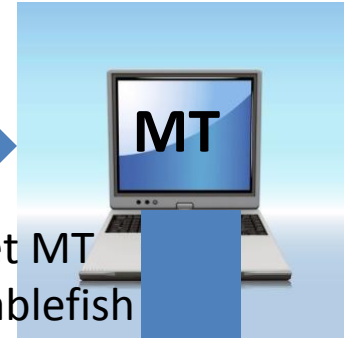**Translation**

**Data Owners**

**End Users**

# Publisher-Centric MT



**End Users**

# User-Centric HT



**Translation**

**End Users**

# User-Centric MT



**End Users**

Free Internet MT
Bablefish
Altavista
Systransoft
Systranet
Google Translate
Free Translation (SDL)
Microsoft Windows Live Translator
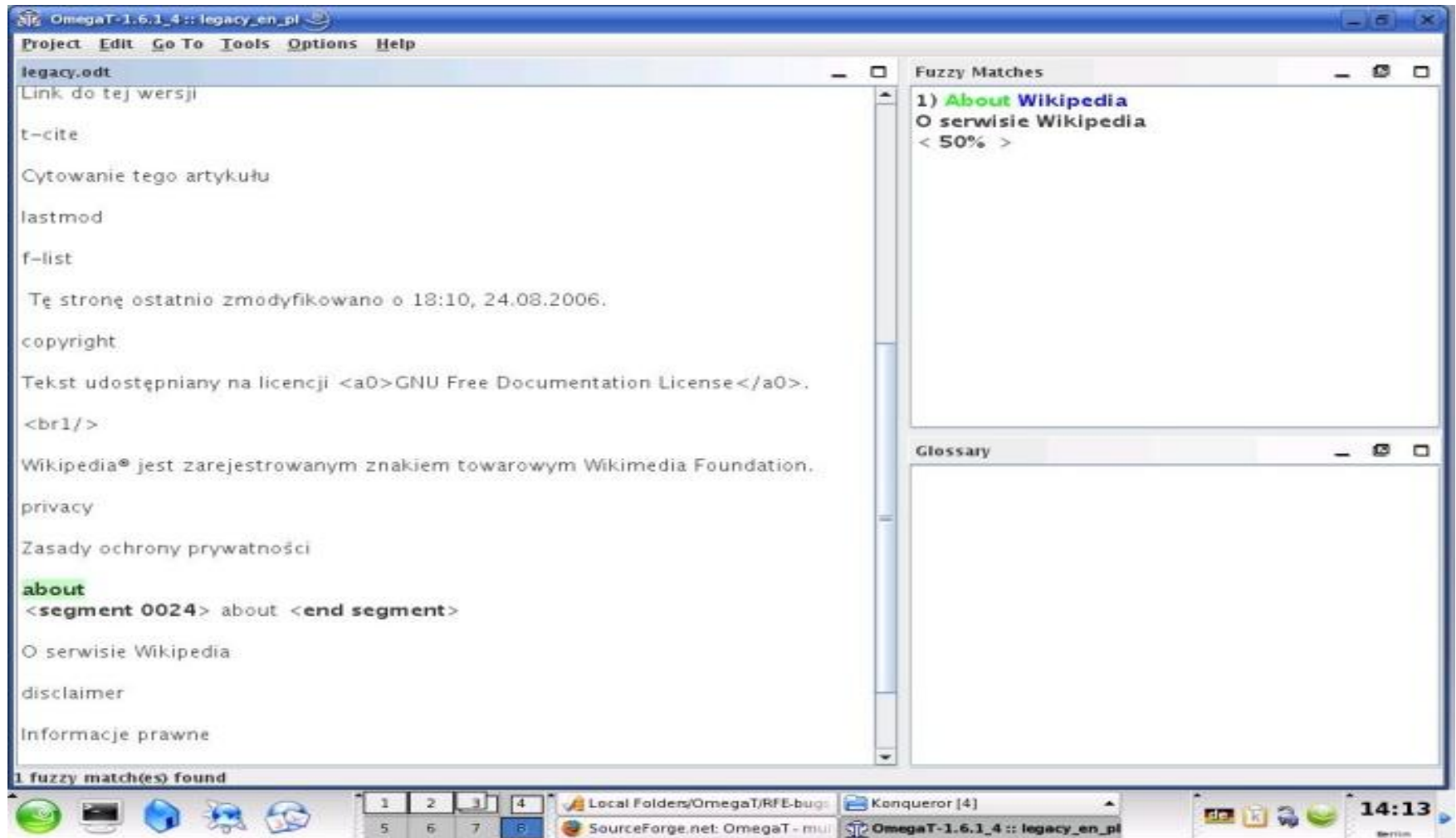Licensed Internet MT
Licensed In-House MT

MT

# Human Translation



Client

Translator

*Dictionaries*
*Wikipedia*
*Chat rooms*

Manager

Editor(s)

*Translation Management*
*Translation Memory*
*Terminology Management*

**Language Service Provider**

End Users

# Translation Memory

- Examples
  - Products:  Across, Bee-Text, Multicorpora, SDL TRADOS, WordFast, LingoTek
  - Open Source:  GlobalSight, Omega-T
- Based on MT research
- Align source and translated texts
- Presents past translations to translator
- Carried to extremes, it provides only text that has not been translated

# GlobalSight Open Source TM

# Terminology Management

- Examples:
  - MultiTerm, Terminotix, Etc.
- Dictionaries of approved terms
  - Content-oriented
  - Sometimes multiple languages
- Often developed by expert translators/ terminologists
- Usually results in clear single right term
  - Potential for MT?

# SDL MultiTerm

# Research of Terms and Contexts

- Search
  - Google, Yahoo, etc.
  - TransSearch
  - Wikipedia
  - Multiple dictionaries
  - Chat rooms
  - Translator site bulletin boards
  - Blogs
- Communicate with
  - Language and language domain experts
  - Authors of text being translated
  - Developers
  - Subject Matter Experts
  - Users (focus groups)
  - Foreign marketing offices
  - Surveys
- Consult
  - Transliteration standards
  - Abbreviation and acronym conventions
- Disambiguate terms
- Document!

# Machine Translation with
# Human Pre-Editing and/or PostEditing



MT

LSP

**Pre-editing**

**Translator**

**End Users**

**Retranslation**

**Post-Editing**

# Pre-Editing

- Examples
  - Authoring systems
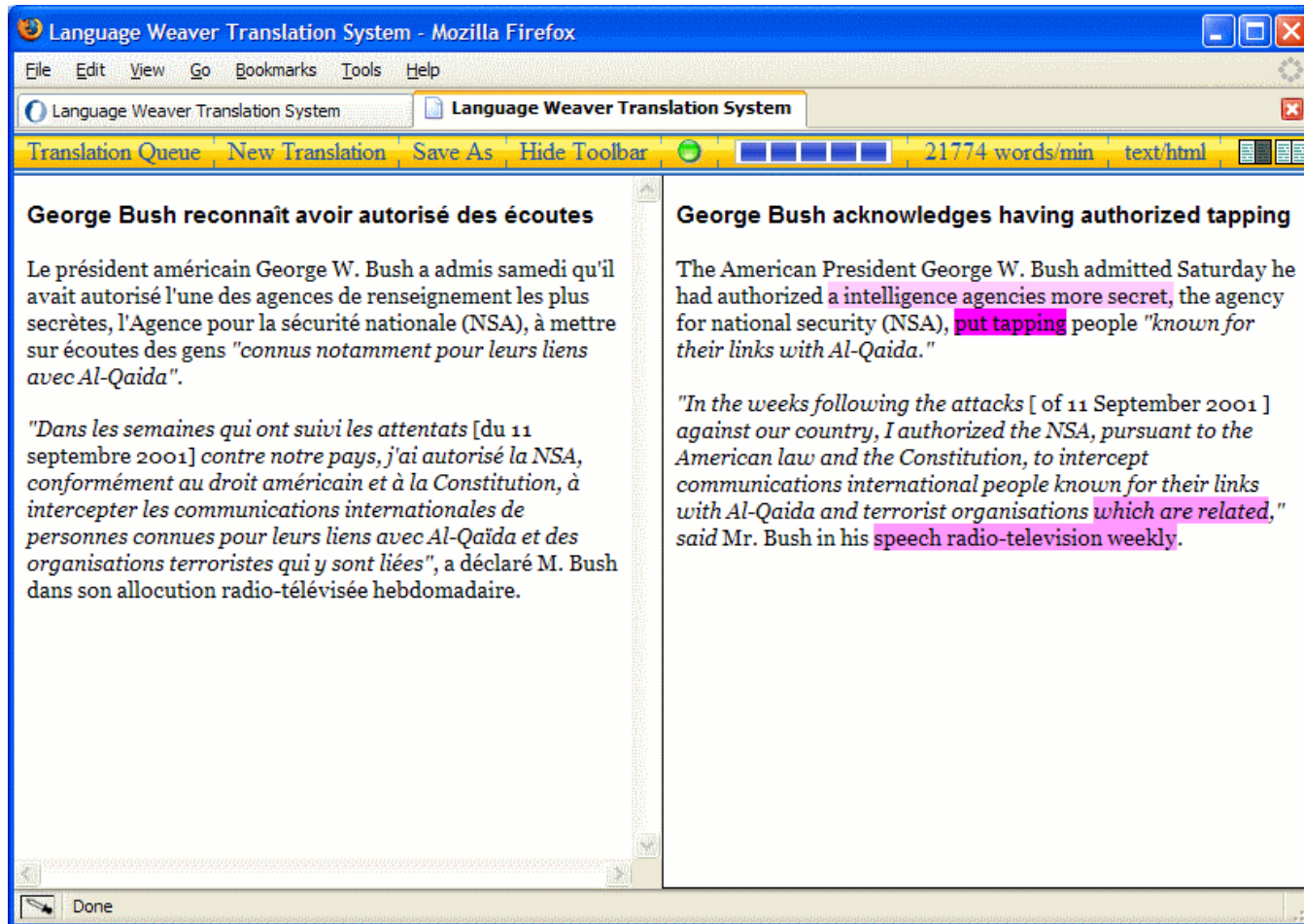    - Xerox, Smart's MaxIT, Acrolynx, and AuthorIT
  - Translatability ratings
    - Bernth and Gdaniec (2001)
- Feedback to authors of text on how to make their writing more machine-translatable
  - Short sentences
  - Full sentences
  - Punctuation
  - Approved vocabulary
- Range from writing tips to controlled language

# Post-Editing

- Examples
- Markup for translators
  - Areas needing work
  - Areas not needing work
- With no markup, translators must read source and target text.  Source text is sometimes not read.
- Interesting editing and QA tools
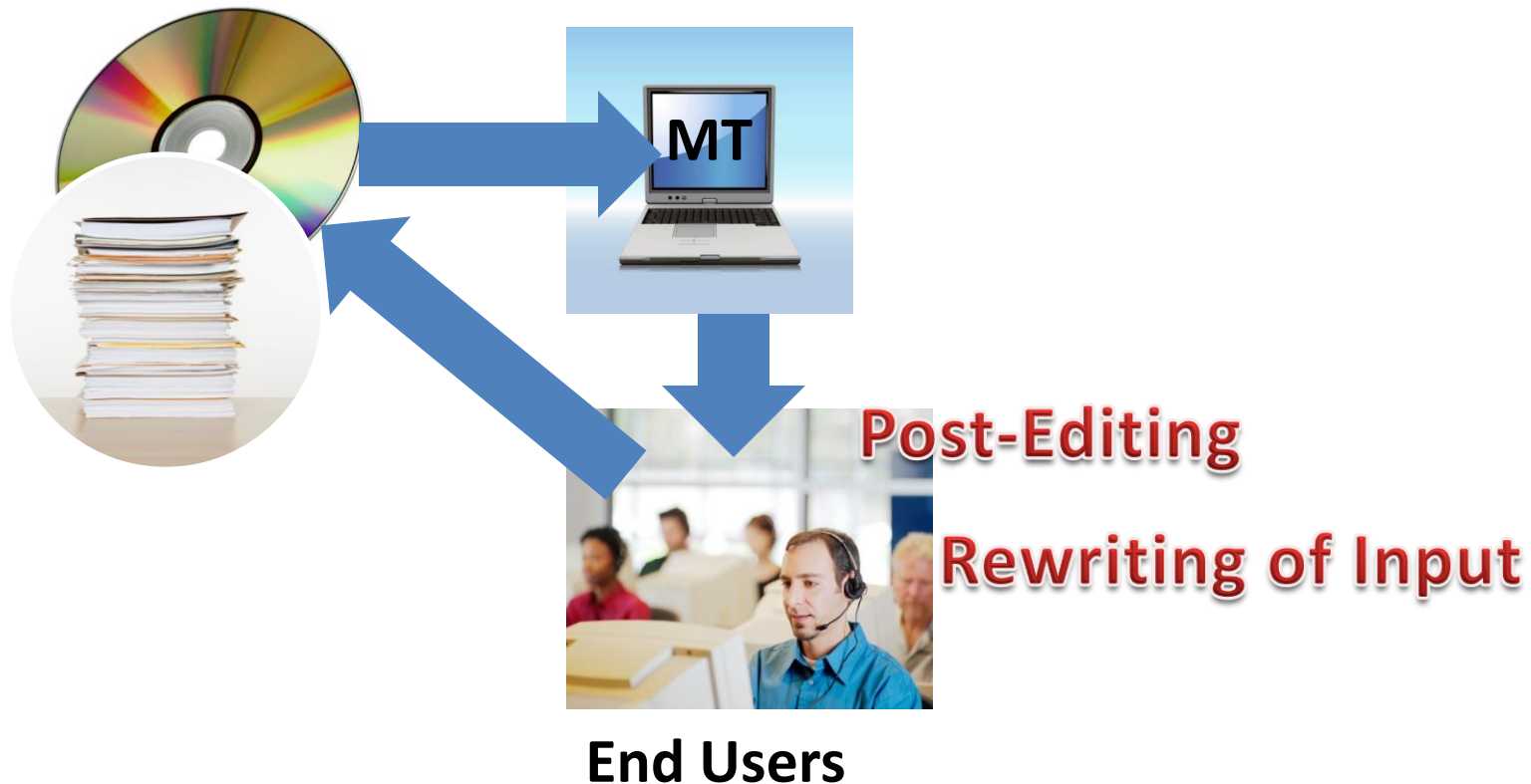  - TransType  (RALI)

# LanguageWeaver Confidence Rating

# Tools for GIL

- Globalization, Internationalization, and Localization (GIL)
  - Globalization, Internationalization, and Localization Technology (GILT)
  - i18n, l10n, g11n
  - Locale
- Tools
  - Strip out text and provide it to translators
  - Help ensure consistency of usage
  - Provide information/means of dealing with differences between countries and regions
    - Language/dialect, voltage info, plug info, font preferences, color and design preferences (Common Locale Data Repository), spelling differences
    - Term differences (terminologies)
      - Useful for MT?

# User-Centric MT with Tools?



**MT**

**Post-Editing**

**Rewriting of Input**

**End Users**

# Gaps

1. Use MT approaches to help HT
2. Use HT approaches to help MT

# How do we better determine and provide term translations?

- What can we do besides provide dictionaries in electronic format?
  - WordNets
  - Translation Memory with context
  - TransSearch (RALI)
  - Morphological Analyzers (right approach?)
- How do we make the right word with the right conjugation appear in the user's translation with the least number of keystrokes?
- How can we share dictionaries and standardize terminology across tools?
  - ISO Lexical Markup Framework, OLIF, etc.

# How do we build better dictionaries and terminologies?

- How do we find new terminology?
  - OOV / NFW searches?
  - Other?
- How do we develop definitions?
  - Can we automate more of the process?
    - Examples, larger term, distinguishing term, etc.
- How do we standardize terms?
  - Frequency is not always the best solution
- How do we provide easy conversion between terms in a translation?
  - Transliterated terms
  - Different terms
- How do we alert users that there may be synonyms?
- How do we disambiguate terms?
  - Frequency is not always the best solution
  - Cycle through by frequency?

# How can we use MT to help translators?

- MITRE Research

- Each source sentence is visually aligned with usually two distinct machine translations of that sentence
  - Sentence boundaries are identified automatically
  - MT output can be copy + pasted into translation/gist authoring pane

- Provides context and overview

- Provides triangulation

# How else can we use MT technologies?

- **TransType**
- EC's Information Society Technologies (IST) Program
- Lapalme, Langlais, Macklovitch, Gandrabur, Nguyen, Nie
- Partners
  - Atos Origin
  - Lehrstuhl für Informatik VI, Computer Science Department, RWTH Aachen - University of Technology
  - Instituto Tecnológico de Informática, Universidad Politécnica de Valencia
  - RALI Laboratory - University of Montreal
  - Celer Soluciones
  - Société Gamma
  - Xerox Research Centre Europe
- Text prediction
  - Saves keystrokes
  - May increase speed and accuracy of translators

# TransType

# TransType:  Options



Do not touch the drum surface.

Touching the drum could reduce the print quality.

WARNING

This warning alerts you to areas of the product where there is the possibility of personal injury.

Si toca el tambor, puede hacer que se reduzca la calidad de impresión.

puede hacer que se reduzca la.
puede hacer que se reduzca la calidad.
puede hacer que se reduzca calidad de impresión.
puede hacer que se reduzca el calidad de impresión.

# Use of MT and TM to Spot-Check Quality

- *Compare source and translated texts to identify problems with*
  - *Omissions*
  - *Numerical expressions*
  - *Source language interference* (Macklovitch)
    - E.g., False cognates
      - Library and librairie
- *Check consistent use of terms*
- *Test translators*

# Use of Tools to Get Bilingual Concordances of Terms

- To
  - Check context
  - Research alternatives
  - Develop terminologies and/or new dictionaries
- Translation Memory
- Search
  - TransSearch (RALI)
  - In house or restricted high-quality search
  - Internet search
- Many opportunities for analysis and development

# TransSearch



Unit handled a 53% increase in the number of terms rendered into French

# TransSearch



*Macklovitch, Lapalme, Gotti, 2009*

# TransSearch
## 20 Most Frequent 2-Word Queries

What could be done with this kind of data?

| Query | Freq. | Query | Freq. |
|---|---|---|---|
| as such | 1195 | as of | 533 |
| over time | 1046 | most importantly | 530 |
| consistent with | 743 | more importantly | 511 |
| in turn | 708 | where appropriate | 499 |
| hard work | 680 | as per | 493 |
| along with | 669 | due diligence | 486 |
| subject to | 655 | build on | 483 |
| based on | 649 | in particular | 478 |
| as required | 609 | focus on | 472 |
| as appropriate | 587 | where possible | 461 |

*Macklovitch, Lapalme, Gotti, 2009*

# Use of Intelligent Templates

- Martin Kay, Xerox PARC
- Restricted options for translators

# How can we get better HT Training?

- MT incorporated into subject-focused training
  - Egan 2009
- Provide automatically-generated lessons on highly relevant subject matter

# How can we improve post-editing?

- Mark the problems for the post-editors
- Provide information on the specific problems
  - Kind of problem
  - Suggested change
  - Linked resources for making change/getting info
- Provide other info
  - Tips on getting better output
- Automate some post-editing (Doyen et al. 2008)

# How can we improve pre-editing?

- Automate some parts
- Make "checkers" less intrusive
- Use source text analysis in additional ways
  - Translatability Ratings (Bernth & Gdaniec 2001)
  - Problem tagging (DeCamp 2009)

# Use of MT, Entity Tagging, and Other Tools

- Examples
  - MITRE research
    - Day et al.
  - CACI Language Workbench

- Help routing

- Provide overview

- Call up resources?

# How Can We Improve Retranslation?

- If a user requests a retranslation of a passage or document, how can we better understand what he/she needs?

- How can we quickly provide it?
  - What additional information can we send without going back to a translator?

# How can we better understand translators? How can we help translators use tools?

- Automated logging
- Recommender systems
- Examples
  - TransType
  - MITRE
- Testing difficult

**CFLEX: Nota6 Translation**

Chart axes: Gist Document Size (characters) vs Time (seconds since start)

Legend:
- keystrokes
- deletes
- pastes
- dictionary lookups

# How Can We Improve MT with HT Practices?

- Footnoting and annotation
- Multiple meanings with examples

# Translator and User Acceptance

- Translators and Post-Editors
- Tolerance depends—and can change
  - Group (demographics, etc.)
  - Tasks and objectives
  - Training and exposure
  - Perception of helpfulness of tools
    - May be influenced by simple issues such as access
- Perception of their customers' requirements
- Difficult to assess productivity gains due to so many factors in human translation

# Conclusions

1. The term "translation" covers many tasks and requirements.

2. There are many means of meeting these requirements with HT, MT, and a combination, depending on needs for accuracy, fluency, etc.

3. There is much room for research and development

    MT to HT

    HT to MT

# At the Technology Showcase

- **Friday, August 28**

- Dictionaries
  - Army Research Labs
  - WordScape
- Translation Memory and/or Tools for Translators
  - Basis Technology (particularly input methods)
  - Canadian Translation Bureau
  - MultiCorpora
  - PAHO
  - ProMT
  - SDL and LanguageWeaver
- Document Triage and Handling
  - CACI
  - Center for Applied MT