2012

AMTA

20Years

The Tenth Biennial Conference of the
Association for Machine Translation in the Americas

# Increasing Localization Efficiency with SYSTRAN Hybrid MT Products

John Paul Barraza

**Systran**

SAN DIEGO, CA
OCTOBER 28- NOVEMBER 1, 2012

This session will cover how to increase localization efficiency with a SYSTRAN desktop product and a server solution. First we will demonstrate how to integrate MT in a localization workflow, interaction with TM matching tools, hands-on MT customization using various tools and dictionaries, and final post-edition using SYSTRAN Premium Translator, a desktop product. We will also walk through the complete cycle of automatic quality improvement using SYSTRAN Training Server, part of the Enterprise Server 7 suite. It covers managing bilingual and monolingual data using Corpus Manager, training hybrid or statistical translation models with Training Manager, and evaluating quality using automatic scoring and side-by-side translation comparison. It also includes other useful tools that automatically extract and validate dictionary entries, and create TMs from unaligned bilingual sentences automatically. Finally, localization efficiency with or without MT integration/customization is compared with the actual cost benefits.

**Presenter**

- John Paul Barraza, Director of Technical Services at Systran.

# SYSTRAN

**Increasing Localization Efficiency with SYSTRAN Hybrid MT Products**

---

## Presenters

— John Paul Barraza, Director of Services
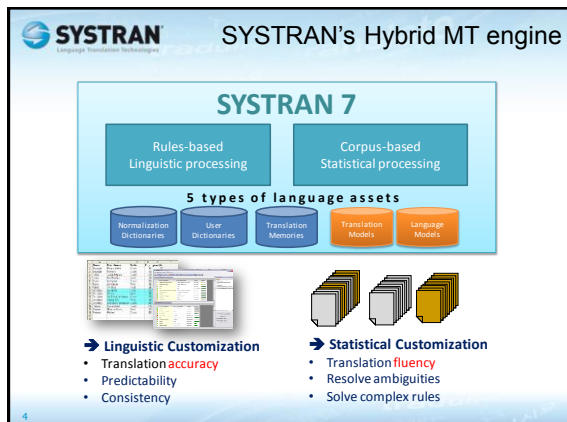
— Philip Staiger, Sr. Technical Trainer

www.systransoft.com

---

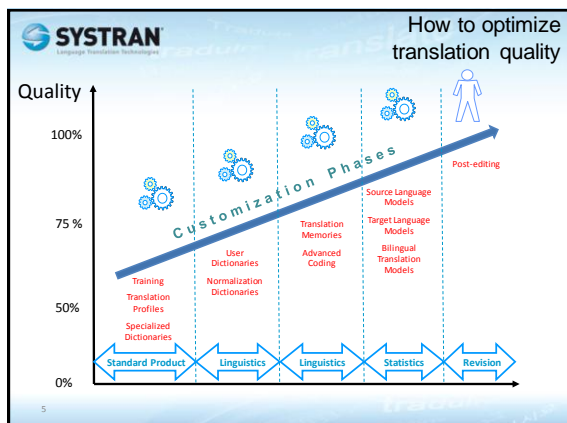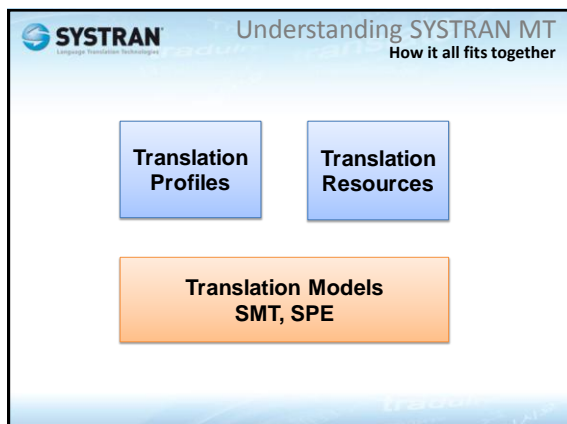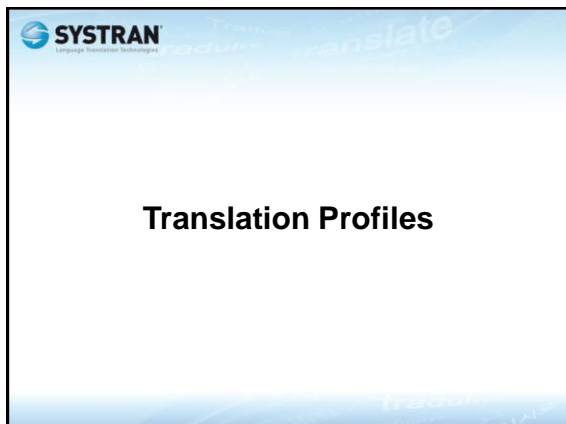## Fast facts

⇒ **Founded in 1968 in CA, now two operating companies**
  ⇒ SYSTRAN SA, in Paris, France (parent company)
  ⇒ SYSTRAN Software, Inc., in San Diego, CA (original comp.)
  ⇒ Publicly traded on NYSE / Euronext (RAN)
⇒ **Business Units**
  ⇒ Software Publishing
  ⇒ Professional Services
⇒ **Machine Translation (MT) Products and Services**
  ⇒ Desktop and Mobile products
  ⇒ Server solutions for the enterprise
  ⇒ Online services (for portals)
⇒ **Advantages**
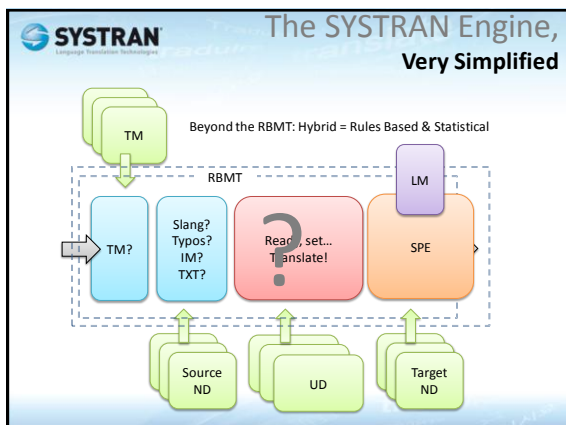  ⇒ Scalable, configurable, customizable
  ⇒ 43+ years experience

## SYSTRAN's Hybrid MT engine

**SYSTRAN 7**

Rules-based
Linguistic processing

Corpus-based
Statistical processing

**5 types of language assets**

Normalization Dictionaries | User Dictionaries | Translation Memories | Translation Models | Language Models

➔ **Linguistic Customization**
- Translation accuracy
- Predictability
- Consistency

➔ **Statistical Customization**
- Translation fluency
- Resolve ambiguities
- Solve complex rules

---

## How to optimize translation quality

Quality

100%

75 %

50%

0%

*Customization Phases*

Post-editing

Source Language Models
Target Language Models
Bilingual Translation Models

Translation Memories
Advanced Coding

User Dictionaries
Normalization Dictionaries

Training
Translation Profiles
Specialized Dictionaries

Standard Product | Linguistics | Linguistics | Statistics | Revision

---

## Understanding SYSTRAN MT
### How it all fits together

**Translation Profiles**

**Translation Resources**

**Translation Models SMT, SPE**

**SYSTRAN**

**Translation Profiles**

---

**SYSTRAN**

The SYSTRAN Engine,
**Very Simplified**

TM

Beyond the RBMT: Hybrid = Rules Based & Statistical

RBMT

LM

TM?

Slang?
Typos?
IM?
TXT?

Ready, set…
Translate!

?

SPE

Source
ND

UD

Target
ND

---

**SYSTRAN**

Translation Profiles:
**Defining Resource Priorities**

Another Example: the Web Interface
➜ **Source & Target Languages?**
➜ **Translation Profile?**

➜ Profiles improve
your Translations:
  – Dictionaries
  – Translation
    Memories
  – Normalization
  – Linguistic
    choices
  – Source
    adaptation
  – Past Translation
    Choices
  – Other resources
    and options

staiger@syst

**SYSTRAN 7**
Enterprise Server

Translation    Dictionary
               Search

Projects: Default.ET...  ▼   **Text Translation**

**Translation**           Choose language source and target. Click on the double arrow to switch sou

**Text Translation**       English  ▼  🔁  Japanese  ▼   Profile **[Default]**

Web Translation

File Translation

RSS Translation

Using a Translation Profile
**Examples**

Inside the SYSTRAN Translation Project Manager ➜

← Inside the Interactive Translator

Inside the Quick file translator (from right-click context menu) ➜



How to Configure a Profile?
**Global Options**

- Resources:
  – Domains & Dictionaries
  – Language Models
  – Translation Choice Files
- Filter Options
  – Formatting
  – Do Not Translate (DNT) Options
- Linguistic Options
  – Source Analysis
  – Country Localization
  – Rendering
  – Imperative
  – Pronouns
  – Style Parameters

SYSTRAN Options



Global Options

**Profile Customization**

## Translation Profiles
### What *can* you Customize?

- Resources
  - Dictionaries
- Filter Options
  - Formatting
- Linguistics Options
  - Language specifics
- Indicators
  - Visual cues

## Types of Resources

- Dictionaries
  - Several types: UD, TM, ND
  - Within a UD: several types of entries
    - Multilingual entries: Source & Target Language
    - SA : Source Adaptation (grammatical categories)
    - DNT : Do Not Translate (acronyms, proper names…)
- Language Models (TM)
- Translation Choice Files (TRC)

## Filter Options
### Do Not Translate (DNT)

color vs. colour

- TM matching: ignore punctuation?
  Ignore Casing? Ignore white spacing?

**Linguistic Options**
Rendering

**Linguistic Options - Rendering**
- ☑ Localize punctuation
- ☑ Localize numbers
- ☑ Add space around Latin characters
- Keep source capitalization
- ☑ Convert Hijri dates
- ☑ Not Found Word transliteration

- Appearance of translated text
  - Punctuation  e.g. 1.00 → 1,00
  - Capitalization

Imperative

**Linguistic Options - Imperative**
Imperative choice - French/German/Italian/Spanish/Portuguese
Infinitive
Imperative choice - Chinese/Japanese/Korean
Polite

| | Example: | Do not lean out! |
|---|---|---|
| German infinitive | → | Nicht hinauslehnen! |
| German imperative | → | Lehnen Sie nicht hinaus! |
| or informal: | → | Lehne nicht hinaus! |
| French infinitive: | → | Ne pas se pencher! |
| Italian (adaptation) | → | E pericoloso sporgersi! |

**Linguistic Options**
Style Parameters

**Linguistic Options - Style Parameters**
- ☐ Separate translation of sentence between quotes
- Document type
  - Default
- ☐ Enhanced technical domain analysis
- ☐ Terminology Translation

**Linguistic Options - Style Parameters**
- ☐ Separate translation of sentence between quotes
- Document type
  - Default
  - Default
  - Abstract
  - List
  - Journalism
  - User Guide
  - Correspondence

- Examples:
  - Scientific papers vs. News Reports vs. Tech Docs vs. Emails/IMs

## Domains:
### Standard vs. User-defined

- Always enabled:
  - Main Dictionary, General Domain
- Also included:
  - Alternative Meaning Dictionary
  - 20 special domains, spread over 5 SYSTRAN System dictionaries:
    - Business Dictionary
    - Industries Dictionary
    - Life Sciences Dictionary
    - Science Dictionary
    - Colloquial Dictionary
- Easily Create your own
  - Manually
  - Wizard (WUD)
- Combine & Prioritize

---

## Tips for Writing

### Source Content Creation for Better Translation Quality

---

## Tips for Writing
### SYSTRAN Web Resources

## Tips:
### Writing Content for Better Translation Quality

- **Be Direct. Write in a Simple, Clear Manner.**
- **Be Concise and To The Point.**
- **Do Not Leave Out Necessary Words.**
- **Beware of Slang and Colloquialisms.**
- **Insert Proper Punctuation & accents.**
- **Check for Accurate Spelling.**
- **Use Articles Whenever Possible.**
- **Consistent Use of Terminology and Abbreviations.**
- **Maintain a Simple Format.**

---

## SYSTRAN
## Customization Wizard

**Leveraging your
Existing Resources**

---

## Customization Wizard:
### Why?

- Existing Terminology ➔ Dictionaries (WUD)
- Target Language Models

## Customization Wizard Output

What are the Resources Created by it?
→ Language Models (non-editable)
→ Wizard User Dictionaries (editable)

Looking at the Results
**From Wizard to Dictionary Manager**

→French and English Terminology: Source & Target Language(s)
→Categories: Nouns, Proper Nouns, Verbs, Adjectives, etc…
→Frequency (e.g. as seen in resource extractions; sortable)
→Default Translation
→Priority (4-Normal, 8-Alternative Meanings only, 9-Disabled)



Using the Results
**Adding 'WUD' to the Profile**



Using the Results
**Language Models**

Source Text → SYSTRAN RBMT → ? → Language Model: A: 15% / B: 0% / C: 85% Choose! → Translated Text

**SYSTRAN**

What Else?
**SYSTRAN Document Aligner**

---

**SYSTRAN**

SYSTRAN
Document Aligner

**Source + Target ➔ Bilingual**

---

**SYSTRAN**

Document Aligner:
**What and Why?**

Document
Aligner

- Have same document in two languages:
  - Source (e.g. EN)
  - Target (e.g. FR)

➔ Produce single document, aligned EN-FR Translation Memory

**Why?**
➔ Use in Translation Profiles
➔ Use in Hybrid Trainings

## SYSTRAN

### Document Aligner
**Step-by-Step**

**Save Your Translation Memory**

→ Save
→ Do not save

**Enter a Name for Your Translation Memory.**

Enter a summary name for the generated Translation Memory.

Aligned - Tips for Writing

---

## SYSTRAN

### What's Next?
**Working with Dictionaries**

---

## SYSTRAN

### SYSTRAN Dictionary Manager

**Terminology Resource and
Translation Memory Management
with SDM**

Normalization Dictionaries

Synonyms (source or target language are the same)

— Initials, Acronyms and Colloquial (chat)
— Common typos



Translation Memories

Import:
- TMX
- Spreadsheet
- XML
- BiText (TAB separated)



Domain Management

## Slide 1

**SYSTRAN** — Domain Management

Examples:
- SYSTRAN Business Dictionary:
  - Economics/Business
  - Legal
  - Political Science
- SYSTRAN Science Dictionary:
  - Computers/Data Processing
  - Electronics
  - Mathematics
  - Mechanical/Engineering
  - Photography/Optics
  - Physics/Atomic Energy
  - Chemistry

➜ Does staking order matter?

## Slide 2

**SYSTRAN** — Profiles & Dictionaries
**Stacking & Priority Rules**

Profile Stacking Order:
➜ Top Dictionaries seen first

Entries can have their own priority!
- 4 (Default)
- 8 (disabled & visible to Alternative Meaning)
- 9 (fully disabled & invisible)

**Resource Extraction, Dictionary Validation?**
➜ Changing entry priority to enable/disable

## Slide 3

**SYSTRAN** — SDM
**Advanced Search Tools**

**Dictionary Filtering**
➜ Entries with Problems?

**Find**
➜ Entries with Space(s)
➜ Separate multi-word entries from single-word entries

## What is the Translation Project Manager?

Why STPM?

- **To manage translation projects**
  - Incremental translation
  - Align source and target sentences
  - Concordance search
  - Alternative meaning display
  - Translation choices save for reuse
  - Feed info to user dictionaries directly from revision tools
- **Multi-file and Multi-format projects**
  - Apply terminology extraction to create term candidates
  - Revise translation with analysis tools
  - Feed directly translation memories from integrated editor
- **Translation comparison**
  - Snapshot creation
  - Sentence complexity and translation accuracy metrics
  - Help to post-edit translations

## STPM for Customization?

Works with SDM & Profiles: new translations will benefit from your choices made in prior translations.

- Crawl: quickly gather extra corpora
- NFW Review: fix Not Found Words (new UDs)
- Also: Mark some NFW as Do Not Translate
- Terminology Review
- Resolve source ambiguities
- Multiple choices? Create Translation Choice files
- Sentence review: keep a Translation Memories

## Translation Project Creation & Management

- File ➜ New
- Single- or Multi-document Project
- Select Language Pair
- Add:
  - Files
  - Folders
  - URL
- Add More:
  - Crawl

Add File
Add Folder
Add URL
Create Folder
Remove from Project
Crawl
Translate
Export Translated Selected Documents
Export Selected Documents as Translation Memory
Project Properties

## SYSTRAN
### 6 Types of Reviews

- NFW Review
  – DNT or Multi lingual entry
  – Export to csv
  – Send to SDM
- Extracted Terms
- Source Ambiguity
- Alternative Meaning
- Sentence
- Translation Memory

TRC File activation in Profile

TRC file: Enable Translation Choices in Profile



Target Ambiguity (Alternative Meanings) Control & Resolution

Alternative meanings found in dictionaries?
- Pick your favorite
- Save as target TRC file



Translation Memory Creation

Sentence Review:
Freeze & Send to the SDM as Translation Memory.

Exporting a Translation as TM



What's Next?
**SMT & Hybrid Translation Models**



**Profiles** ✓   **Resources** ✓

**Translation Models
SMT, SPE**

Hybrid Technology Overview

RBMT, SMT, SPE



Before Hybrid Training:
Rules-Based MT Architecture



Statistical MT Architecture

**SYSTRAN**   **Hybrid MT Resources**



**SYSTRAN**   Training Server:
**Basic Workflow**

- 2 things needed in Training Server
  - Corpus Manager: assemble your materials
  - Training Manager: train your models
- SYSTRAN Enterprise Server SES7: Online platform for complete MT customization
  - Manage & analyze your bilingual / monolingual data (corpora)
  - Train and optimize MT systems for the best quality
  - Evaluate & Compare the customized systems in detail



**SYSTRAN**   Corpus Manager

- Repository used for training processes
  - Monolingual corpora
  - Bilingual/Multilingual corpora
- Upload many supported file types
  - TMX, Text, MS Office, RTF, PDF, HTML,… +**ZIP(!)**
- Database (MySQL) storing TUs (Translation Units)
- Corpus Search to verify
- File system view
  - Hierarchical folder/sub folder structure
  - Virtual Files
  - Partitioned Files (Training, Testing, Tuning)

## SYSTRAN

### Training Manager
### Training new MT models

- Training Manager uses data in Corpus Manager
- Manage various training processes
  - Baseline Translation to Establish Reference Scores of RBMT
  - Bilingual Terminology Extraction
  - Dictionary Validation
  - Hybrid Training for Hybrid MT with SPE
  - Statistical Training for Statistical MT (SMT)
- Task Management
  - Launcher, Monitoring Current Activity, Statistics
  - Automatic quality scores
  - Task Comparator (comparing 2 training runs)

## SYSTRAN

SYSTRAN Training Manager:

### Hybrid MT Resources

- Create additional MT resources to boost Hybrid MT quality
  - Resource Extraction
    - Create UD from bilingual corpus by BTE (Bilingual Terminology Extraction) technology
    - Create Source Adaptation model from monolingual corpus
  - Dictionary Validation
    Validate UD entries against bilingual corpus
  - Document Alignment
    Create TMX by aligning source and target TUs from not-aligned bilingual documents

## SYSTRAN

### Complete Training Cycle:
### What 5 Steps are Involved?

1. Gather a Corpus: bi-lingual phrases, Translation Memories, good quality, additional mono lingual sentences
2. Prepare the Corpus for Training (Load & Partition)
   - 3 Subsets: Training, Tuning, Testing
3. Baseline Evaluation: Reference Scores of the current RBMT?
   - If you care to compare the progress of your models
4. Hybrid Training (or Statistical): several iterations
   - Extract Resources, get new UDs
   - Validate old & new UDs
   - Run Hybrid Training
   - Compare scores: Tweak parameters & Re-Train with various settings
   - Check quality: Compare translated phrases from Testing corpus
5. Publish the best Translation Model, Create new Profile

**SYSTRAN**

Translation Workflow
&
Training Workflow

---

**SYSTRAN**           Translation Workflow:
**SYSTRAN inside the MT pipeline**



- Translation Memory Matching, Fuzzy Scores
- Good enough? Send to human post editing
- Different? Get SYSTRAN involved to pre-translate it; then send to human post editor

---

**SYSTRAN**           Training Workflow:
**Goal of a Training Cycle**

- Generate Improved Translation Model
- Publish it to Translation Server
- Use it in new Translation Profiles

Benefit of this Exercise:
- Improve the quality of the pre-translated output
- Reduce effort & time of human post editor
- Reduce production cost

The Ultimate Goal:
**Smarter Profiles, Better Translation Models**

- Enhance the Rules Based System, e.g. with new UDs
- Enhance the Statistical Post Editing (SPE) with bilingual Translation Models
- Unsupported language? Create purely statistical model when there's no rules engine
- Smarter Technology, Better Translation



Corpus Manager

**Step-by-step**



After you login
The Full Interface

Translation Server

Training Server

**SYSTRAN**

Working within a Project

➔ Select a Project
➔ Or create a new one

---

**SYSTRAN**

Reminders:
What we're about to do

➔ **Gather the Corpus**
  – Bilingual files (TMX, BiText…)
    • Optionally additional monolingual files
  – Partitioning: Getting it ready for Training
    • 3 parts: Training corpus, Tuning corpus, testing corpus

➔ **Train a new Translation System**
  – Measure (Score) the current system: Baseline as Reference
  – Train, several iterations as needed:
    • Run Training (Hybrid, or purely Statistical)
    • Compare the Scores
    • Extract Resources, Validate User Dictionaries
    • Tweak parameters & re-Tune / re-Train as Needed
  ➔Publish the new resources

---

**SYSTRAN**

Corpus Manager
step-by-step

---

Reminder: BiText Format
**Easily create your own TM**

```
#ENCODING=UTF-8
#TM                          ← this is a Translation Memory!
#EN<TAB>ES                   ← use a 'TAB' to separate EN and  ES
First sentence.<TAB>Primera frase.
Got a second phrase.<TAB>Tengo una Segunda frase.
Three is a charm.<TAB>Tres son muy buenos.
```



Corpus Explorer



Corpus Manager
Uploading a File

How to upload
a file into the
Corpus Manager:

- Select File
- Pick Source Language
- Select Filter type
- Directory
- Domains & Comments

→ Got Many files?  Upload a Zip archive of the files





Got (too) many files?

→ Wrap them into a single "virtual" file

→ Even if from different file types

→ Treat it as a single entity for Training purposes

1. Select the files in the explorer view
2. Click 'Create a New Virtual File'
3. File Name: give your virtual file a name
4. Path: set Directory location for the virtual file
5. Add an optional comment
6. Indicate if you want all of their content (100%) or just a (randomly selected) subset
7. Go!

**SYSTRAN**

Virtual Files
**Another Use**

Another Common Scenario:

– Is your corpus too big for just a quick test?

➔ Use a Virtual File with a percentage of less than 100%, for example 5% of the original complete corpus.

➔ Example: Corpus of 1 million TUs down to 50K TUs

➔ Test with small Subset, Experiment

➔ Final Training with the entire corpus

**SYSTRAN**

# Partitioning

Preparing the Corpus
for Training

**SYSTRAN**

Partitioning: Why?
It's all about Scoring and Comparing

- Create 3 subsets: Training (big), Tuning, Testing
- How we train the system:
  – Use the bulk (90%-95%) of the material for training purposes (but not all: keep some for iterative fine tuning, some for testing, i.e. scoring)
  – Set aside a small amount of the Corpus for testing purposes (no overlap: unknown to the trained content)
  – Set aside another small amount of the material for tuning the training process
  – Example 1: 90% for training, 5% for testing, 5% for tuning
  – Example 2: 95,000 training, 3000 tuning, 2000 testing

SYSTRAN

**Partitioning**
## Are there Alternatives?

Use handy Web interface to partition, or:
- Create 3 files and upload them into a folder:
  - training_corpus.txt
  - testing_corpus.txt
  - tuning_corpus.txt
- Have additional monolingual corpora?
  - mono_EN.txt
  - mono_FR.txt
  - mono_DE.txt
  - etc…

SYSTRAN

**Partitioning**
## If you create your own

Avoid Overlaps – don't cheat yourself
- Segments in one corpus should not exist in any of the other two corpora:
  - Training, tuning, testing
  - You don't want to run tests on what's already been seen in the training

**SYSTRAN**

Ready to use the corpus
for training purposes

Next: Training Manager

**SYSTRAN**

Training Manager

Overview of Tasks,
Concepts, Process

**SYSTRAN**

**Task Launcher**

Baseline Evaluation

**Hybrid Training**

Statistical Training

Resource Extraction

Dictionary Validation

Document Alignment

W            rver?

## SYSTRAN Training Scenario: Sample Workflow

**SYSTRAN**

1) **Build your Corpus**
   – Optional Tool: Document Alignment

   **Corpus Manager**

2) **Baseline Evaluation**
   – Reference Scores?

   **Training Manager**

3) **Train a model**
   – SPE, SMT
4) **Optional: Resource Extraction & Dictionary Validation**
   – Improve the RBMT: Disambiguation, Source Adaptation, new Terminology
5) **Repeat & Publish**
   – Translation Models, User Dictionaries, Language Models…

➔ Practically speaking: train multiple systems with different training options and select the best system based on the resulting score.

---

**SYSTRAN**

## Training Manager:

### Baseline Evaluation

---

**SYSTRAN**

## Baseline Task
### Reference Score Evaluation

- **Goal**: generate **Reference Scores** for comparisons
- **Why**: Used to evaluate the results of subsequent Trainings.
- **How**: Use Bilingual Testing corpus + Translation Profile
- **Also**: Compare Translated output vs. Test Corpus Reference Translation

| S | Bilingual Test Corpus | T |

RBMT

? → WER GTM BLEU TER

**Next Steps?**

- We've run a Baseline eval (RBMT Reference)
  - Testing corpus
- We have two more Corpora yet unused:
  - Training corpus
  - Tuning corpus
- Ready to train… but which way?
  - Consider & Choose?  Maybe both
    - Hybrid (Rules & SPE) translation?
    - Purely Statistical translation?



**Two Types of Training**
Hybrid vs. Statistical Training

Task Launcher

Baseline Evaluation

**Hybrid Training**

Statistical Training

Resource Extraction

Dictionary Validation

Document Alignment



**Hybrid vs. Statistical Training**

Improved the BLEU score from 22.9 to 65.6 in under 50 minutes of training.
(config: SYSTRAN servers running in a VM w/CentOS 5, 8GB RAM, hosted by a 16 GB Win7 Pro system with 1st generation i7)

Bilingual Terminology
Extraction (BTE)

**Automatic Creation of Dictionaries**



SES 7
Training Server

Resource Extraction

### The Hybrid Engine Process
#### Where BTE can help

**SYSTRAN**

SYSTRAN HYBRID MT ENGINE

S → [ Source Analysis | RBMT | Target Disambi-guation ] → SPE → T

**Source Analysis:** Disambiguation of term categories (adjective, noun, verb,...) can improve and speed up the RBMT process.

**Use of Terminology:** single- and multi-word vocabulary placed in new dictionaries that will help the translation quality.

**Target Disambiguation:** Mono-lingual Language Models to further improve the fluency.

---

### Purpose of Resource Extraction
#### What to expect from BTE

**SYSTRAN**

- Multi-word user dictionary  ← perhaps the most valuable
- Single-word user dictionary
- Source adaptation dictionary
- Candidates (another UD)
  – Possibly more words, but you need to translate
- Target Language Model

➔ Add these resources to a Profile after publishing them back to the Translation server

---

### After Running a BTE
#### What to do with it

**SYSTRAN**

- In a perfect world:
  – UDs are ready for use as is
  – Publish to Translation Server
  – Add new UDs to a new or existing Profile
  – Retrain: get better Translation model
- In the real world:
  – Some items in the extracted UDs may not help.
    • Filter them (run Dictionary Validation task)
  – Publish & Use the filtered versions
  – Export to csv for human review

Resource Extraction
Advanced & Expert Options



BTE Example
BLEU Score Improvement

➔ With large sets, automatic resource extraction can often yield improvements of 3-8 points on BLEU score, sometimes even more.
➔ BTE is probably the most significant improvement that can be observed with relatively little manual work: no coding.
➔ Additional improvements observed when followed by dictionary validation (discussed later)



Another BTE Example
13 TMX's w/Monolingual Corpus

**SYSTRAN**

Training Manager:

**Dictionary Validation**

---

**SYSTRAN**

SES 7
Training Server

Dictionary Validation

Task Launcher
- Baseline Evaluation
- Hybrid Training
- Statistical Training
- Resource Extraction
- Dictionary Validation
- Document Alignment

---

**SYSTRAN**

Do you Really Need DV?
Can you trust your dictionary entries?

**Problem:** Some entries don't help the translation, make it worse
- Old / Outdated Entries / Out of Domain or Context
- Typos

**Solution:** Find which entries make it Better or Worse if enabled
- Run Dictionary Validation Task
- DV takes time and lots of memory
- For each entry in the dictionary:
    1. Finds all occurrences of the entry throughout the corpus
    2. Translates these occurrences: with and without the entry in use
    3. Metric: How many times was it Better, Same, Worse?
        ➔ If most often it's Better with the entry enabled: enable it
        ➔ If most often it's Worse with it enabled: disable it
        ➔ If most often it's the same, don't change the entry's status

**SYSTRAN** — What to expect from DV
Filtered dictionaries

**SYSTRAN 7** — Enterprise Server

→ DV task completed, New UDs Ready, what now?
- Publish the Task's result, see the filtered UD
- same name, plus 'filtered'
- Add filtered UDs to the Profile(s) & re-Run Baseline test to validate improvement
- Retrain Hybrid model, using new Profile with new UDs



**SYSTRAN** — Launching a DV task
Step-by-step

Launch Dictionary Validation Task



**SYSTRAN** — DV
Recommendations

- Repeat DV for each dictionary that you'll want to use in future profiles and/or trainings, such as:
  - Your pre-existing dictionaries, especially old ones
  - New MultiTerm dictionaries that were never validated
  - Multiword from BTE
  - Singleword from BTE
- Resource Extraction can use a lot of resources. Use subset of corpus if too large
  - 100,000 – 200,000 entries in bilingual corpus is enough

Putting it all together



Complete Cycle (1/4)
Bilingual Corpus, English-Spanish

➔ Uploaded Bilingual EN-ES Corpus
➔ Partitioned: training, tuning & testing



Complete Cycle (2/4)
Baseline tests & Resource Extraction

Baseline with Default profile:

Baseline with MultiTerm user dictionary in profile:

After Resource Extractions: Singleword and Multiword Dictionaries:

Complete Cycle (3/4)
Entity Recognition Rules & DV



Complete Cycle (4/4)
Hybrid Trainings

More than Doubled (!) the BLEU score

Original: BLEU 26.87  TER 53.05  ➔ Final Trained Model:  BLEU: 58.52  TER 26.20

## Conclusion

This completes our exploration of how to increase localization efficiency with SYSTRAN Hybrid MT Products.
Here are some of the topics we've covered:

- Translation Profiles
- User Dictionaries (UD)
- Normalization Dictionaries (ND)
- Translation Memories (TM)
- Linguistic Resources
- Resource Extractions (BTE)
- Dictionary Validation (DV)
- Hybrid Training, Statistical Training