# Air Force Research Laboratory

## Reversing the Palladius Mapping of Chinese Names in Russian Text

**31 October 2012**

**Dr. Katherine M. Young**
**N-Space Analysis, LLC**

**Human Trust & Interaction Branch (711HPW/RHXS)**
**711th Human Performance Wing**
**Air Force Research Laboratory**
**Wright-Patterson AFB**

*Integrity ★ Service ★ Excellence*

# **Acknowledgments**

Co-authors:

Jeremy Gwinnup

Joshua Reinhart

SRA International, Inc.

Thanks to:

Dr. Ray Slyh, AFRL

Dr. Tim Anderson, AFRL

Mr. Bill McIntyre, USAF

# Overview

- The Problem
- Workflow
- Interface Design
    - Human-in-the-loop
    - Automatically-generated options
    - Saving a dictionary
- Applying the Dictionary to New Documents
- Future Possibilities

# The Problem

When using Systran to support Russian to English translation:

- Systran translates known words, and sounds out unknown words
- Chinese names follow a separate Palladius sound mapping

| | |
|---|---|
| Meaning | "British expert Phillip Clark" |
| Cyrillic | Британский эксперт Филлип Кларк |
| Systran | The British expert Of fillip Clark |

| | |
|---|---|
| Meaning | "Zhai Zhigang was the youngest of six children" |
| Cyrillic | Чжай Чжиган был младшим из шести детей |
| Systran | Chzhay Of chzhigan was young of six children |

# Palladius Sound Mappings

Pinyin    Cyrillic    typical sound

| Pinyin | Cyrillic | typical sound | |
|---|---|---|---|
| r | ж | [ž] | Palladius maps initial zh to чж |
| ch | ч | [č] | |
| zh | чж | [č ž] | |
| | | | |
| n__ | н__ | [n] | Palladius maps word-final ng to н |
| __n | __нь | [n'] | |
| __ng | __н | [n] | |
| | | | |
| xi | си | [ɕi] | Palladius maps x to c |
| sa | ca | [sa] | |

# Reverse Palladius (RevP) Algorithm

**1. User selects candidate words**

> <u>Чжай Чжиган</u> был младшим из шести детей

**RevP**

**2. Syllabification**

> #Чжай#  #чжи#ган# был младшим из шести детей

(C)(G)V(N)

C=consonant
G=glide {i, u}
N=coda{n,ng,i,u,r}

**3. Reverse Palladius mapping**

> <u>Zhai zhigang</u> был младшим из шести детей

<u>RevP Rules</u>
чж → zh
н# → ng

**4. Systran**

> <u>Zhai zhigang</u> was young of six children

# Workflow -- Original

The user compiles a list of Chinese/Russian syllable correspondences and uses this list to hand-edit the Systran errors.

Cyrillic
Чжиган  →  Systran  →  errors
                        chzhigan

List
ган > gang
Чжи > Zhi   →   User   →   hand-corrected output
                           Zhigang

- The user has to look back at the original spelling of the name.

- Each instance in the text has to be separately corrected.

- These corrections must be repeated in each new document.

- The user selects a word and the program provides the Reverse Palladius form
- The new mappings are saved in a dictionary for future use, either with Systran or as a stand-alone process.

Cyrillic
Чжиган

Systran

better automated output
in future documents
Zhigang

User

RevP

dictionary
Чжиган
> Zhigang

corrected words in
current document
Zhigang

# Interface Design

- We use a human-in-the-loop to identify the candidate words

- We provide alternative mappings

- We store the user's choices in a dictionary

- The dictionary can be applied to new documents within Systran

- The dictionary can be applied independently

We don't want to apply the Reverse Palladius mapping to Russian words,
   or to words borrowed from other languages

Cyrillic                                                    Лисов

Traditional Sound Mapping                    Lisov        << correct form

Reverse Palladius Mapping                    Lisuow

Cyrillic                                                    Малинди

Traditional Sound Mapping                    Malindi      << correct form

Reverse Palladius Mapping                    Malingdi

If the name occurs as the object of a verb or preposition,
it may appear with Russian case endings.

| Cyrillic | ReversePalladius | pinyin | |
|----------|------------------|--------|---|
| Бомин | Boming | Boming | |
| Бомина | Bominga | Boming | (genitive) |
| Боминóм | Bomingom | Boming | (instrumental) |

We optionally apply a stemming program to remove inflections
before the sound mapping:

| meaning: | " ... together with Liu Boming" |
|----------|----------------------------------|
| Cyrillic: | ... вместе   с   Лю Боминóм |
| Reverse Palladius: | Liu Bomingom |
| Stemming + Reverse Palladius: | Liu Boming |

# Human-in-the-loop:
# Identifying Russian Inflection

We cannot apply the stemmer in all instances.

For example, the Chinese name <u>Baohua</u> ends in /a/, which is a potential Russian inflectional ending.

Stemming of <u>Баохуа</u> creates an error:

| Cyrillic | Stemmed | RevP | Correct pinyin |
|----------|---------|------|----------------|
| Баоху<span style="color:red">а</span> | Баоху | Baohu | Baohua |

The RevP program generates the various possible forms, and we rely on the user to select the correct alternative.

# Interface: User Selection



The user has selected the name, Бомин.

The program presents the user with the various options for the selected word. The first option is the Reverse Palladius Mapping.

# Interface: Replacing the Word



**ReversePalladiusUI**

File   Edit   Dictionary   Options   Help

Экипаж «Шэньчжоу-7»
состоял из трех хантянъюаней (космонавтов):
№1 - Чжай Чжиган,
№2 - Лю Boming,
№3 - Цзин Хайпэн.

Лю Бомин тоже родился в бедной семье, вторым из шести детей.

Семьи Чжай Чжигана и Лю Бомина и сегодня живут в г. Цицикар.

**Replace All**

Would you like to replace all occurrences with 'Boming'?

[ Yes ]   [ No ]

The selected form replaces the original, and the program offers to replace all instances of this word.

# Interface: Replacing the Word

Here, the user has selected the word Бомина.
Now we need the stemmed form, which is the second option.

# The RevP options:

| P | Bominga |
|---|---------|
| Ps | Boming |
| Cp | Bomina |
| S (Systran) | Systran |
| Manual Entry | Manual Entry |
| book | No Dict. Match |

| Symbol | Source | Details |
|--------|--------|---------|
| P | Palladius | apply reverse Palladius mapping |
| Ps | stemmed Palladius | remove Russian inflectional endings, then apply reverse Palladius mapping |
| Cp | Cyrillic phonetics | apply Cyrillic phonetic mapping |
| S | Systran translation | use the Systran translation for this word (if you have Systran enabled) |
| hand | Manual Entry | type in the correct spelling |
| book | Dictionary Match | use the previous dictionary choice for this word (if any) |

Systran translation:  a preview of how Systran would translate this word or phrase

- Systran 7 SOAP API – use code generated from WSDL to allow RevP to send highlighted text to a Systran server for immediate translation

# Interface:  Manual Entry



Manual Entry:  If none of the provided forms are correct,
there is an option to type in the correct form.

# Interface: Manual Entry

# Interface:  Manual Entry

# Saving the Dictionary



In this example, the user has processed the names Liu and Boming.  We select the "View Dictionary" option…

# Saving the Dictionary



… and see the two entries that RevP has created.

# Saving the Dictionary



These entries can be stored to a file by using the menu option, File/Save Custom Dictionary.

The dictionary is saved in a Systran-readable format.

# Applying the Dictionary
# to New Documents Using RevP

We can use the RevP program to apply a saved dictionary to a new document.

In this example, we have created a dictionary with the names of three cosmonauts.



cosmonauts.rpd - Notepad

```
#ENCODING=UTF-8
#SUMMARY=REVP-Systran/cosmonauts.rpd
#MULTI
#RU        EN
Хайпэн     Haipeng
Лю         Liu
Чжай       Zhai
Цзин       Jing
Бомин      Boming
Чжиган     Zhigang
```

# Applying the Dictionary to New Documents Using RevP



We start the RevP program and open a new document.

# Applying the Dictionary
# to New Documents Using RevP



We select the menu option, File/Open Custom
Dictionary, and load our cosmonaut dictionary.

# Applying the Dictionary to New Documents Using RevP



The program pre-translates all the matching names in the document.

We can also compile our cosmonaut dictionary as a Systran user dictionary.



```
cosmonauts.rpd - Notepad

File  Edit  Format  View  Help

#ENCODING=UTF-8
#SUMMARY=REVP-Systran/cosmonauts.rpd
#MULTI
#RU        EN
Хайпэн     Haipeng
Лю         Liu
Чжай       Zhai
Цзин       Jing
Бомин      Boming
Чжиган     Zhigang
```

# Applying the Dictionary
# to New Documents Using Systran



## Select the dictionary tab, select dictionary compilation, and compile the dictionary file.

Select the translation tab and enter the Russian text (or upload a file).

# Applying the Dictionary
# to New Documents Using Systran



Before applying our user dictionary, we get transliterated versions of the Chinese names:
chzhay chzhigan, tszin khaypen, lyu bomin

# Applying the Dictionary
# to New Documents Using Systran



We click "Select your dictionaries" and check the RevP dictionary.

# Applying the Dictionary
# to New Documents Using Systran



Now we get the correct forms:  Zhai Zhigang, Jing Haipeng, and Liu Boming.

- Integrate a named entity tagger to help select the Chinese names

- Include an option for Tongyong pinyin

    (used in Taiwan)

- Tag the corrected names as do-not-translate items going into Systran

- Adapt the program to other languages, to assist in information retrieval

# Questions?

Katherine M. Young, PhD
N-Space Analysis, LLC
(937) 475-3637
nspaceanalysis@earthlink.net