## MTTT – Machine Translation Training Tool A Tool to Teach MT, Evaluation and Post-editing

P. Bouillon (FTI), P. Estrella (FaMAF), R. Lafuente (FaMAF), M. Lemos (UBP), S. Girletti (FTI)

- Collaboration: FaMAF@UNC, UBP, Córdoba, Argentina & TIM/FTI@UniGe, Geneva, Switzerland
- **Goal**: integrate complete MT+PE workflow into one tool
- **Scope**: for pedagogical use
- Advantages: open-source tool, GUI, desktop and web versions, user guide provided
- Current state: on-going development
- Further development: access to intermediate results, tmx support, extensive testing





#### **Desktop version**

#### Web version

#### **Corpus Preparation**



ous Preparation	Model training	Machin	ne Translation	Evalu	ation	Post-Editing	Diff - State
La	nguages			1	Fransla	tion Model —	
Source	de	:	Source text	mme	ntary-	v8.de-en.de.txl	:
Target	en	:)	Target text	mme	entary-	v8.de-en.en.tx	t
	Results			_	Langu	age Model —	
/home/paula/Do	cuments/linux-6	4bi 🔒	Text (target	lang)	htary-	v8.de-en.en.tx	t
en <	ter/tokenizer.per	140			Se	ttings	
/home/paula/Do commentary-v8. /home/paula/Do	wnloads/news- de-en.en.txt > ocuments/test/in	DU	Output Dire	ctory	ula/D	ocuments/tes	:
t.tok.en	,,			Start	corpus	preprocessing	1

8 🖻 🗊 Machi	ine Tra	nslation Training	g Tool				_	
Corpus Prepar	ation	Model training	Machine Translation	Evaluation	Post-Editi	ng Dif	f - Stats	
	Evaluation							
MT (Hyp)	ents/	/final-test/test-er	n-de.txt.translated		WER	PER	MTER	
Reference	Reference ula/Documents/final-test/test-en-de.txt				BLEU			
	Start Evaluation							
WER: 1.06 PER: 1.06 HTER: 0.9: BLEU: 97.22	Vert L readouri WER: 1.06 (Execution errors:) PER: 1.06 (Execution errors:) HTER: 0.32 (Execution errors:) BLEU: 97.22 (Execution errors:)							

a 🖨 🖲 Machine Translation Training Tool						
Corpus Preparation Model training Machine	Translation	Evaluation	Post-Editing	Diff - Stats		
Search						
Differences	Differences Stats					
And if you so, those who are really worried that, gold might indeed mouidmine the m reliable hedge worried	d about D nost e A	Deletions: 0, Insertions: 0, Replacements: 3 Times edited: 2 , Approx. time spent in minutes: 1				

#### GitHub

roxana-lafuente/MTTT PaulaEstrella/MTTT-PyQT



Corpus Preparation	Taining	Machine Translation	Evaluation	Post-Editing	Differences
	Un	redited NT			Edited MT
SAN SAN FRANCISCO den Wert von Gold zu 1	<mark>0 -</mark> Es war no führen.	och nie leicht, ein rationa	íles Gespräch i	berSAN <mark>FRAN</mark> von Gold zu	Es war noch nie leicht, ein rationales Gespräch über den Wert 1 Tühren.
in letzter Zeit allerding: letzten Jahrzehnt um ü	s ist dies schu iber 300 Proz	vieriger denn je, ist doct rent angestiegen.	ı der Goldpreis	im In letzter Zei letzten Jahra	it allerdings ist dies schwieriger denn ja, ist doch der Goldpreis im 2 zehnt um über 300 Prozent angestegen.
Erst letzten Dezember Roubini Kommentare, Marktstimmung hinterh hinviesen.	verfassten m in denen sie r ragten und se	eine Kollegen Martin Fe nulig die vorhertschend ehr überlegt auf die Risik	idstein und Nou e optimistische zen des Goldes	riel Erst letzten I Roubini Kon Marktstimm, hinwiesen.	Dezember verfassten meine Kollegen Martin Feldstein und Nouriel 3 mmentane, in donen sie multig die vorhenschende optimistische ung hinterfragten und sehr überlegt auf die Risken des Goldes Mandear wie in Hinterregeingen

#### miguelemosreverte/MTTT\_Web







# SCATE Smart Computer-Aided Translation Environment Year 3 / 4

#### http://www.ccl.kuleuven.be/scate

Vincent Vandeghinste, Tom Vanallemeersch, Liesbeth Augustinus, Frank Van Eynde, Joris Pelemans, Lyan Verwimp, Patrick Wambacq, Geert Heyman, Marie-Francine Moens, Iulianna van der Lek-Ciudin, Frieda Steurs, Ayla Rigouts Terryn, Els Lefever, Arda Tezcan, Lieve Macken, Sven Coppers, Jan Van den Bergh, Kris Luyten, Karin Coninx



## Poster Topics

- Semantic-based alignment
- Syntactic concordancing
- Translation with probabilistic synchronous tree substitution grammar
- A neural network architecture for detecting grammatical errors in SMT
- Domain adaptation for ASR Post-Editing
- A deep classifier for bilingual lexicon induction
- Intelligible and in-context recommendations for translation environments



Your single-source partner for corporate product communication

20th Annual Conference of the European Association for Machine Translation (EAMT)

## TM & MT – a happy couple ...or how to calculate the potential benefit

Elsy Stolze, STAR Language Technology & Solutions

### What added value does MT provide?

- ▲ For customers with...
  - an existing translation process
  - extensive Translation Memory and terminology
- ▲ We offer: 3-phase proof of concept
  - Initial analysis
  - Productive pilot phase
  - Productive analysis



Prof. Dr. rer. nat. Alexander Ferrein Gary Evans B.Sc.

foundcat.fh-aachen.de | cat.fh-aachen.de www.fh-aachen.de Part of the Fellowship -Programm 2016 | www.stifterverband.org/digital-lehrfellows



Ministerium für Innovation, Wissenschaft und Forschung des Landes Nordrhein-Westfalen



#### 50 Students translated 10 scientific articles in 1 week homework Peer reviewed.



Hi Guys,

This computer aided translation tool is a very good idea because it helps many people when they are searching for high-quality translations, which you can't find so much on the Internet.

Thanks to this project you can help other students and improve your English skill at the same time, so I did a couple of translations, I hope it will help somebody :)

Comp	onents	Languages	Informati	on Glossaries	History
Tools	✓ Sh	аге 🔻			
Compo	onent 🗸		Translated	✓ Words √	Тга
"Black Fik	out": Zwi	schen	62	2.9% 5	1.8% A
An	Reviewed To be revi	: 100.0% ewed: 0 of 18 segr	nents		
Au Lei	FH unte	rstützt "March for	Science"	The FH supports "Ma	rch for Scienc
Da				✓ Rev	view

Prof. Dr. rer. nat. Alexander Ferrein Gary Evans B.Sc.

foundcat.fh-aachen.de | cat.fh-aachen.de www.fh-aachen.de Part of the Fellowship -Programm 2016 | www.stifterverband.org/digital-lehrfellows



Ministerium für Innovation, Wissenschaft und Forschung des Landes Nordrhein-Westfalen







An open-source neural machine translation system.

English Français 简体中文 한국어 日本語 Русский العربية

> @harvardnlp @systran

OpenNMT is a industrial-strength, open-source (MIT) neural machine translation system utilizing the Torch/PyTorch toolkit.



- Features include:
  - Simple general-purpose interface, requires only src/tgt files.
  - Speed/memory optimizations for **multi-GPU** training.
  - A dependency-free **C++ translator** for model deployment.
  - Latest research features to improve translation performance.
  - Extensions to other sequence generation tasks:
    - Text summarization and
    - Image-to-text generation.
  - Active **community** with academic and industrial contributions.
  - Pretrained **models** available for several language pairs.



An open-source neural machine translation system.

English Français 简体中文 한국어 日本語 Русский العربية

> @harvardnlp @systran

 State-of-the-art results in WMT 2017 (English-German news task)

#### Ongoing research projects

- Efficient data sampling strategies for NN training.
- Better initialization for fast optimization convergence.
- Probabilistic line search approach for SGD. (https://arxiv.org/abs/1703.10034)
- Multi-encoders for neural machine translation (https://arxiv.org/abs/1601.00710)
- Linear mapping (bridge) between encoder-decoder layers.
- Modeling coverage for NMT. (https://arxiv.org/pdf/1601.04811)
- Domain control for NMT. (https://arxiv.org/abs/1612.06140)
- Hyper-specialization techniques for NMT















	Interagency Factsheet on refugee and migrant children and UASC in Europe 2016
- 8	🛃 Télécharger 🛃 Détails   1.59 MB
	Téléchargements: 9.675

60







UNION EUROPEA Fondo Social Europeo





# Is there a case for the use of MT in the Third Social Sector?

@celiaricoperez Universidad Europea



## **OpenNMT Toolset (OMTS)**



- Streamline the process of creating workable NMT models
- Help users choose the best model by evaluating OpenNMT output
- Integrate neural machine translation into enterprise localization process
- Enable users to try the latest machine translation technology with least effort

Minimal Dependency Translation: a Framework for CAT for Under-Resourced Languages

 For under-resourced languages (URLs), lack of written material

processing languages of the global south

afaan oromoo

guarani

ghichwa

k'iche

- Translation: a partial solution
- But insufficient resources for CAT for URLs
- MDT: phrase-based RBMT for URLs
- Assumed resources
  - SL: tokeniser, POS tagger, morphological analyser
  - TL: morphological generator
- Project goals
  - Open-source tools for implementing MDT
  - Implementations for Spanish–Guarani, English–Amharic



#### TraMOOC project

Aims:

- a translation platform offering reliable translation services for 11 languages
- target languages include weakly supported and/or morphologically rich languages
- various types of educational MOOC texts presentations, assignments, video lecture subtitles, forum blog texts

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

#### TraMOOC project

Recent results:

- translation systems based on PBMT and NMT
- comparative analysis of two approaches:
  - automatic evaluation
  - human rankings of adequacy and fluency
  - error annotation
  - technical and temporal post-editing effort
- NMT: better ranking, better fluency, lower error count however, several aspects are not clear

# SkyCode MT – a translation system using deep syntactic and semantic analysis

Luchezar Jackov, SkyCode Ltd.

Universal dictionary based on Princeton WordNet

 The members of analysis and laviage synthesis are based on

The morphological analysis and lexical synthesis are based on lexicalizations bound to PWN synsets

- Language-independant semantic knowledge base The knowledge base is built around dependency relations on PWN synsets
- Analysis is split from synthesis and analysis and synthesis rules are shared among similar languages

#### Compact data representation

The total space needed for 42 translation directions is approx. 150 MB

# SkyCode MT – a translation system using deep syntactic and semantic analysis

Luchezar Jackov, SkyCode Ltd.

#### Portable implementation

The system is implemented in C++ which makes it portable across various platforms

 Currently available for English, German, French, Italian, Spanish, Bulgarian and Turkish

 Easy adding of new languages allowing translation from and to all of the languages in the system

• The rule-based nature of the system makes it easy to implement under-resourced languages

### ModernMT: A New Open-Source Translation Platform for the Translation Industry











Real time adaptation to content

No more training

or set-up time



Designed to scale for data and users



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 645487.

#### www.modernmt.eu

#### ModernMT: A New Open-Source Translation Platform for the Translation Industry

- Currently: classical phrase-based approach
- New data can be added at any time
- Phrase table entries are created on the fly ...
   ... by phrase extraction from a selection of parallel text most similar to the input
- Interpolation of several language models chosen from a large pool of domain / task-specific LMs
- On Github: http://github.com/modernmt

# amadeus

#### GLOBAL DISTRIBUTION SYSTEM

#### Challenges



Provide custom machine translation of Fare Quote Notes for 5,000+ travel agencies located in Russia, Armenia and Belarus

#### Business requirements

- Understandable ("good enough") translation of Fare Quote Notes (FQN) from English into Russian with focus on frequent terms and abbreviations
- High performance and reliability even in case of the large number of requests

# Integration in the Amadeus system's interface



## PROMT Cloud API

Cloud solution with programming interface (API), ready to process large number of translation requests

### Customization

Pre- Processing	Additional Dictionaries	Translation Memory	<ul> <li>End-users</li> </ul>	s feedba	ck	
Special	~1,200 terms	professional	Usage stat	istics		
take into account FQN format and structure	<b>~20,000 names</b> of airline companies and airport codes	translation of the most frequent sentences	4 million words per week	6 reque per w	k ests /eek	~700 words in every request

#### Implementation results

Translation of an entire FQN

Additional functions:

PROMT translation takes into

Integration in the Amadeus system:

Translation of current page of an FQN

account linguistic features of FQN

Translation of special terms on hover

**PROMT** 



# <u>Automatic Building of Machine Tran</u>slation (2013-2016)

#### **Selected final results**

- Goal rapid and cost-effective MT for new languages
- ✓ Work academia ↔ industry secondments
- Case study Croatian as a new official EU language
- Results publicly available tools and datasets

### Data acquisition

#### **Machine Translation**

















# Appraise on Azure

- What is this about?
  - A cloud-based, multi-purpose evaluation framework
  - Used for WMT evaluation campaign but supports arbitrary annotation tasks!
- What has changed?
  - Latest frameworks: Python 3, Django 1.11
  - Now runs on Azure
    - either inside Ubuntu VM, or
    - natively on Windows Server



## • Why should I care?

Appraise on Azure allows focus on evaluation, not install/build an eval system



# Azure for Research

Cloud computing awards Big data, <u>real</u> HPC, machine learning, deep learning (GPUs), IoT, and more....

Deadline: 15 June 2017

(and every two months)

- Email:
  - azurerfp@microsoft.com







WWW.azure4research.com

#### HimL: Health in my Language

Type: Horizon 2020 Innovation Action Duration: February 2015 to January 2018 Website: www.himl.eu









Charles University, Prague





#### HimL: Aims

We will make public health information available to consumers in their own language by:

- Deploying MT which is:
  - Adapted to the domain
  - Semantics- and discourse-aware preserves meaning
  - Can translate into languages with rich morphology found in central and eastern Europe
- Integrating the MT into the content management workflow of the two users (Cochrane and NHS 24)
- Carefully managing user expectations
- Comprehensively evaluating user satisfaction and impact of the new functionality



A translation-based approach to the learning of the morphology of an under-resourced language Tewodros Abebe, Addis Ababa University Michael Gasser, Indiana University

- Most languages have no computational morphological resources and insufficient data for creating them
- Most under-resourced languages are closely related to better-resourced languages
- Objective: learn the morphology of an under-resourced (target) language using translation from a related well-resourced (source) language
- Wolaytta (source) and Gofa (target)
  - Omotic family of southwestern Ethiopia
  - Wolaytta: best studied, Gofa: very few resources



## **Architecture and Processing**





Your single-source partner for corporate product communication

20th Annual Conference of the European Association for Machine Translation (EAMT)

> MT in real-world practice: Challenges & solutions at Swiss Federal Railways

Nadira Hofmann, STAR Language Technology & Solutions

### Introducing MT at the Swiss Federal Railways...

- Preconditions & motivation
- ▲ Approaches & requirements
- ▲ Deployment
  - > Workflow system
  - SBB Translate

## Babeldr vs Google Translate: A User Study at Geneva University Hospitals

P. Bouillon, J. Gerlach, H. Spechbach, N. Tsourakis and S. Halimi

## Aim: compare two translation tools for diagnosis in emergency settings

DI UNIVERSITÉ HUG Lineratares		Google	Sign in
BabelDr		Translate	Turn off instant translation
depuis combien de jours avez-vous mal au ventre ?	منذ كم يوم تنشغرُ بالم في البطن؟	French Arabic English Detect language -	Arabic French English + Translate
		depuis combien de jours avez-vous mal au × ventre	كم يوما هل لديك ألم في المعدة
	Patient(e) Domaine Voix du système	<b>4) Ů III ·</b> 47/5000	☆ 🗍 ♠) 🗳 🖉 Suggest an edit
Enregistrer Traduire	arabe (HOMME)		kam yawmaan hal ladayk 'alam fi almueadd
rechercher ici	D: depuis combien de jours avez-vous mal au ventre ?		
depuis combien de jours avez-vous mal au ventre ? depuis combien d'heures avez-vous mal au ventre ? avez-vous mal au ventre depuis longtemps ? la douleur au ventre est-elle récente ? avez-vous mal au ventre [ <i>ily_durée</i> : depuis longtemps]? avez-vous mal au ventre [ <i>ily_durée</i> : depuis longtemps]? <i>depuis_durée</i> : depuis longtemps <i>depuis_durée</i> : depuis longtemps		Google Translate for Business. Translator Toolkit	Website Translator Global Market Finder
dopuio, durée: dopuis plus d'une hours Projet Bi	abelDr	About Google Translate Community Mobile G+	About Google Privacy & Terms Help 🗾 Send feedback

ISSCO/



# Results

- o 18 diagnoses
  - o 9 French doctors
  - 2 standardized Arabic-speaking patients
- GT is not precise enough for the anamnesis, even if some doctors could reach a correct conclusion

BabelDr

Was not perceived as less robust and limited

Allowed more correct diagnoses



# Dissecting Human Pre-Editing toward Better Use of Off-the-Shelf Machine Translation Systems

Rei Miyata (Nagoya University) Atsushi Fujita (NICT)

## **Objectives**

1) Investigate the capability of the pre-editing strategy

- A human-in-the-loop protocol to collect pre-edit instances
- Japanese-to-English translation tasks on 4 datasets
- 2) Provide an **overview of possible edit operations** 
  - A typology of edit operations

# Source text (ST)



# MT output

Birth registrations submitted to the public office of the municipality where you live.

Please submit notification of birth to the public office of the municipality where you live.



## **Protocol for Collecting Pre-Edit Instances**

Human editors incrementally edit source texts (STs) relying on their introspection, so that improved MT quality is achieved. (Miyata et al. 2015)

New features: (1) Record every minimal edit, (2) Allow reversion of past edits



4 domains (400 sentences) Ja-En SMT Ja editor ↓ 12,287 pre-edit instances

## ✓ More than 85% STs achieved satisfactory MT quality

✓ English-translatable Japanese STs are also Chinese- and Korean-translatable







#### **Mtradumàtica**

Adrià Martín-Mor (UAB) Sergio Ortiz-Rojas (Prompsit) Gökhan Doğru (UAB)



#### Statistical Machine Translation Customisation for Translators

- Free (GPL) Moses-based platform
- User-friendly web interface
- Installable in private servers
- Understand how the components of the system work together
- m.tradumatica.net

#### Prepare your SMT system in 5 simple steps



Fast initial training (10x faster than Moses) Seamless assimilation of new user data (while decoding) Real-time learning (feature computation) On-the-fly adaptation (from context) Elastic architecture (from single to multi-node) Ready-to-install package (docker or binaries)

