



October 28 - November 1, 2016



The Twelfth Conference of  
The Association for Machine Translation  
in the Americas

<http://www.amtaweb.org/amta-2016-in-austin-tx>

---

November 1, 2016  
**ModernMT**

Marcello Federico  
Marco Trombetti

# ModernMT Tutorial

Marcello Federico - FBK, Italy  
Marco Trombetti - MateCat, Italy  
Davide Caroselli - Translated, Italy

AMTA - 1 November 2017, Austin, Texas

ModernMT

Next Generation  
Machine Translation



## Outline

- Introduction + (Marcello, 40 min)
- Development + (Davide, 20 min)
- *Break*
- QA and discussion (all, 20 min)
- Hands-on (Marco & Davide, 70 min)

ModernMT

Next Generation  
Machine Translation

# ModernMT Introduction

Marcello Federico - FBK, Italy

AMTA - 1 November 2017, Austin, Texas

ModernMT Next Generation  
Machine Translation



## Translators' pains with MT

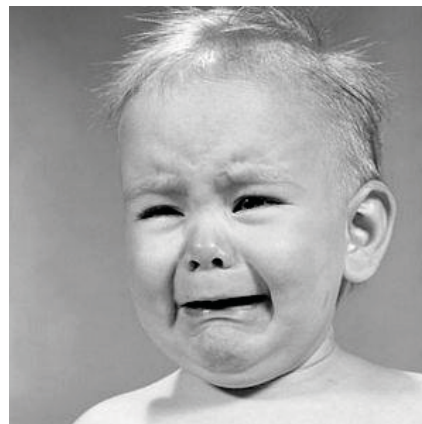
output is often **poor** or **contextually wrong**



ModernMT Next Generation  
Machine Translation

## LSP engineers don't laugh either

**cumbersome setup of MT**  
**lack of training data**  
**online MT is too generic**



## Setting up MT for CAT today

- (1) select TMs
- (2) collect extra data
- (3) train and evaluate engine
- (4) doesn't work? back to (2)
- (5) apply MT on documents
- (6) create fake TM with MT
- (7) import TMs in CAT tool
- (8) start translating
- (9) adapt engine to new data
- (10) new project? back to (1)

# The Modern MT way

- (1) connect your CAT with a **plugin**
- (2) drag & drop your **private TMs**
- (3) start translating!



**ModernMT** Next Generation Machine Translation

# Enhance translator's experience



Adaptation from TMs  
Learning from corrections

**ModernMT** Next Generation Machine Translation

# Modern MT in a nutshell

**fast** training  
**manages context**  
**learns** from users  
**scales** with data and users



Modern**MT** Next Generation Machine Translation

# Team

Business



Research



Modern**MT** Next Generation Machine Translation

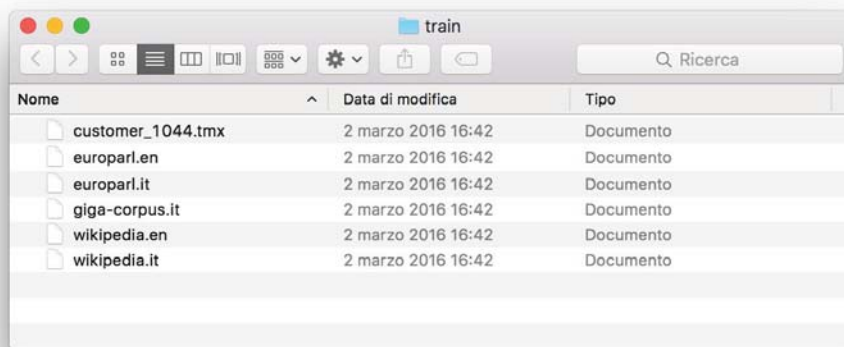
# Acknowledgment

program: H2020  
 type: innovation action  
 funding: 3M €  
 duration: 2015-2017  
 grant: 645487



**ModernMT** Next Generation  
 Machine Translation

# Fast and easy training



```
> mmt create en it path/to/data
```

**ModernMT** Next Generation  
 Machine Translation

# Prototype - Fast training

Training takes **30s** for a  
**1M word TM**

**MMT** is **12x time faster**  
than std Moses



# Context aware translation

TEXT 1

**We're going out.**  
**party**

TRANSLATION

**Nous sortons.**  
**fête**



# Context aware translation

TEXT 1

**We're going out.**  
**party**

TEXT 2

**We approved the law.**  
**party**

TRANSLATION

**Nous sortons.**  
**fête**

TRANSLATION

**Nous avons approuvé la loi.**  
**parti**

# Context aware translation

SENTENCE

party

CONTEXT

**We are going out.**

CONTEXT

**We approved the law**

TRANSLATION

**fête**

TRANSLATION

**parti**

# REST API

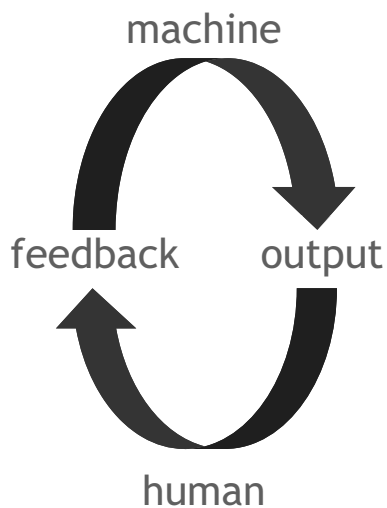
GET /translate?q=party&context=We+approved+the+law

```

{
  "translation": "parti",
  "context": [
    {
      "id": "europarl",
      "score": 0.10343984
    }, ...
  ]
}
    
```

**> mmt start**

# learning from users



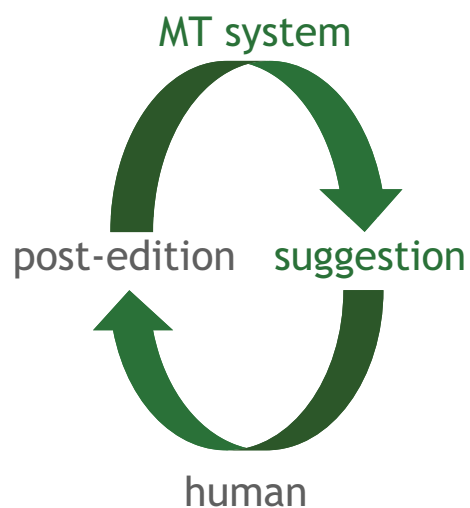
## machine

- > **supports** the user
- > **improves** over time

## human

- > keeps her workflow
- > is an **unaware trainer**

## incremental learning



### goal

- > **correct** recurring MT errors
- > keep **system** up-to-date

### features

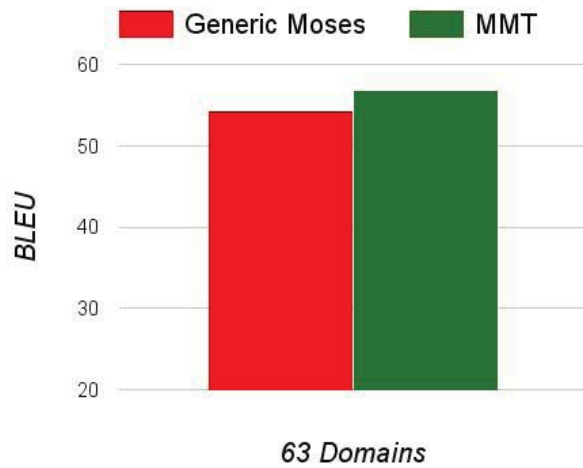
- > **adapt** to user/domain
- > **real-time** processing

## Enhanced text processing

- **More** supported languages
- **Faster** processing
- **Simpler** to use
- **Tags** and **XML** management
- Localization of **expressions**
- **TM cleaning**



# Prototype (March 2016)

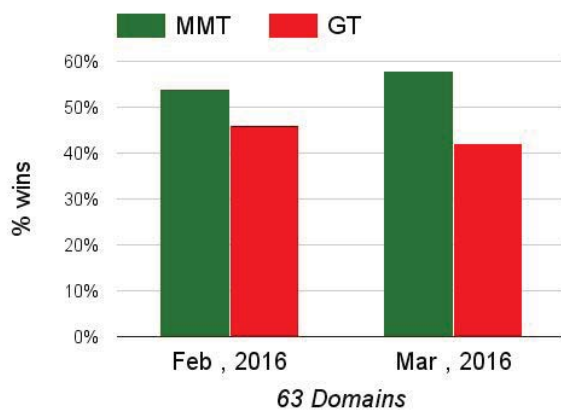


**Benchmark 1.1**  
- No XML tags

**MMT outperformed generic Moses by >2 BLEU points**  
**12x faster training**

**Speed: 1.7sec/segm**

# Prototype (March 2016)

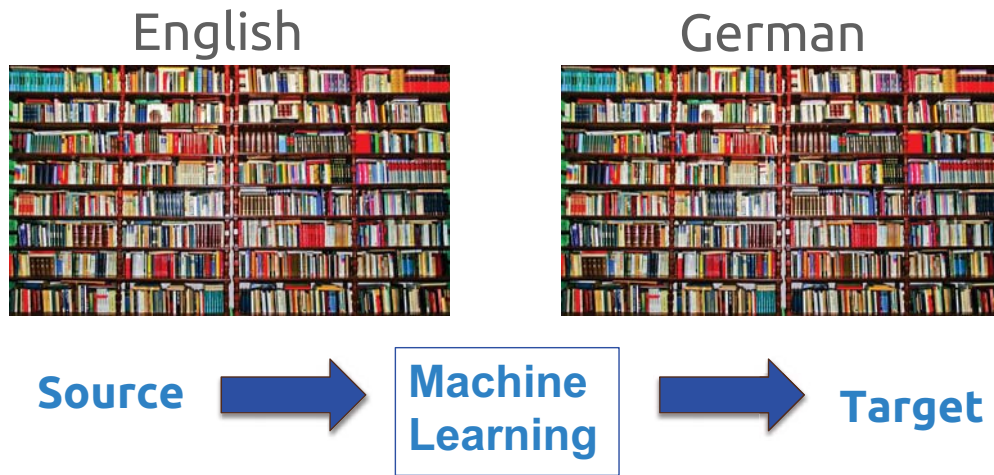


**Benchmark 1.1**

**A/B testing vs GT:**  
~ 300 rnd segments  
~ 3 judges

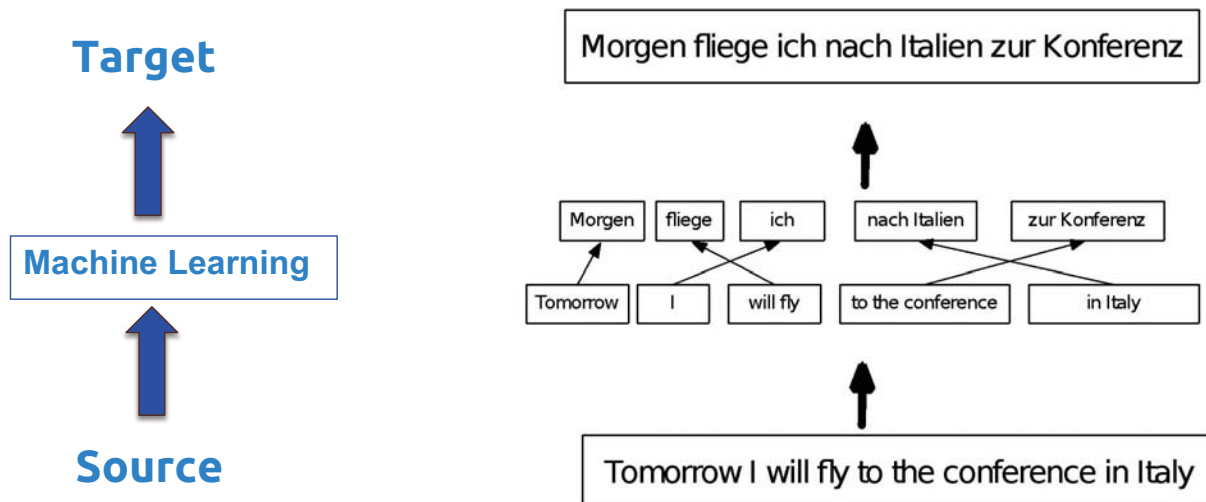
**Distance doubled, from 8% to 16%!**

# Core Technology



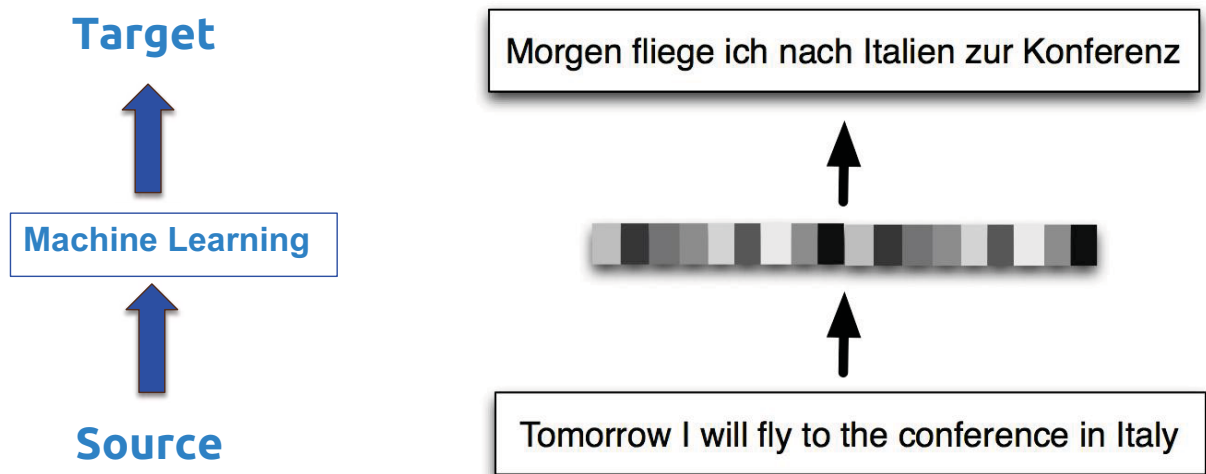
ModernMT Next Generation Machine Translation

# Today: Phrase-based MT



ModernMT Next Generation Machine Translation

# Tomorrow: Neural MT (2017)



**ModernMT** Next Generation Machine Translation

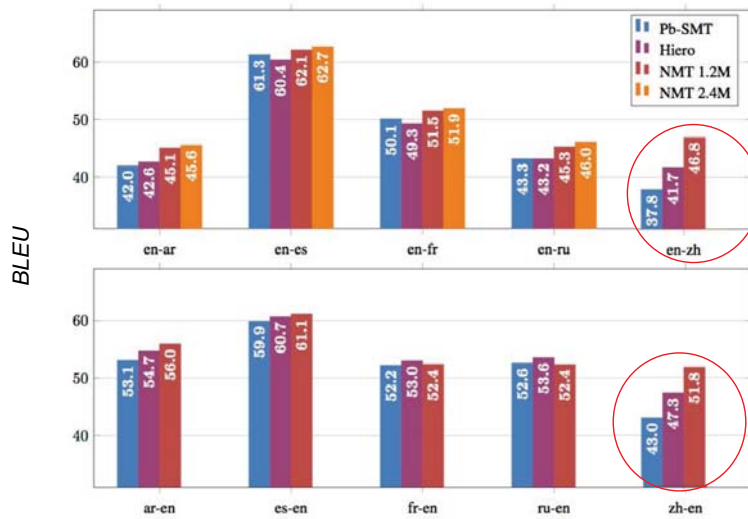
# Phrase-based versus Neural MT

- ❑ PBMT has dominated for about 20 years
- ❑ NMT started to outperform PBMT in 2015
- ❑ NMT is now taking over PBMT in research

	PBMT	NMT
Computational cost	0.80\$/hour (CPU-16)	0.70\$/hour (GPU-1)
Training time (100Mw)	few hours	1 week
Adaptation (1Mw)	< 1 min	hours
Translation speed	455 w/s	409 w/s
Translation quality		+ 5-25% BLEU

**ModernMT** Next Generation Machine Translation

# Phrase-based versus Neural MT



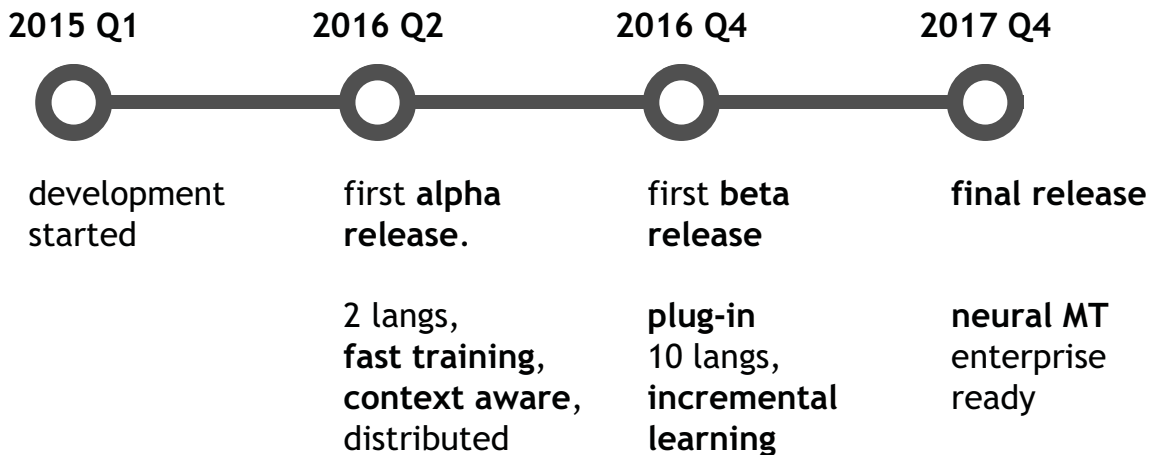
Task: United Nations Parallel Corpus v1.0.

Source:

M. Junczys-Dowmunt, T. Dwojak, H. Hoang. *Is Neural Machine Translation Ready for Deployment? A case Study on 30 Translation Direction*, arXiv. 11 Oct 2016.

**ModernMT** Next Generation Machine Translation

# Roadmap



**ModernMT** Next Generation Machine Translation

## Did I mention that MMT will be free?

LGPL/Apache licences  
new core technology  
no licensing



[github.com/ModernMT/MMT](https://github.com/ModernMT/MMT)

---

Modern**MT** Next Generation  
Machine Translation

## Conclusion

# Modern MT

big data, context aware, enterprise

---

Modern**MT** Next Generation  
Machine Translation



# Thank You

**Project website:**  
**[www.ModernMT.eu](http://www.ModernMT.eu)**



**[github.com/ModernMT/MMT](https://github.com/ModernMT/MMT)**

## Outline

- Introduction + (Marcello, 40 min)
- Development + (Davide, 20 min)
- *Break*
- QA and discussion (all, 20 min)
- Hands-on (Marco & Davide, 70 min)

# Modern MT Development

Davide Caroselli - Translated, Italy

1 November 2016 - Austin, Texas

ModernMT Next Generation  
Machine Translation



## The Modern MT way

- (1) connect your CAT with a **plugin**
- (2) drag & drop your **private** TMs
- (3) start translating!



ModernMT Next Generation  
Machine Translation

## Modern MT in a nutshell

**fast** training  
manages **context**  
**learns** from users  
scales with data and users



---

Modern**MT** Next Generation  
Machine Translation

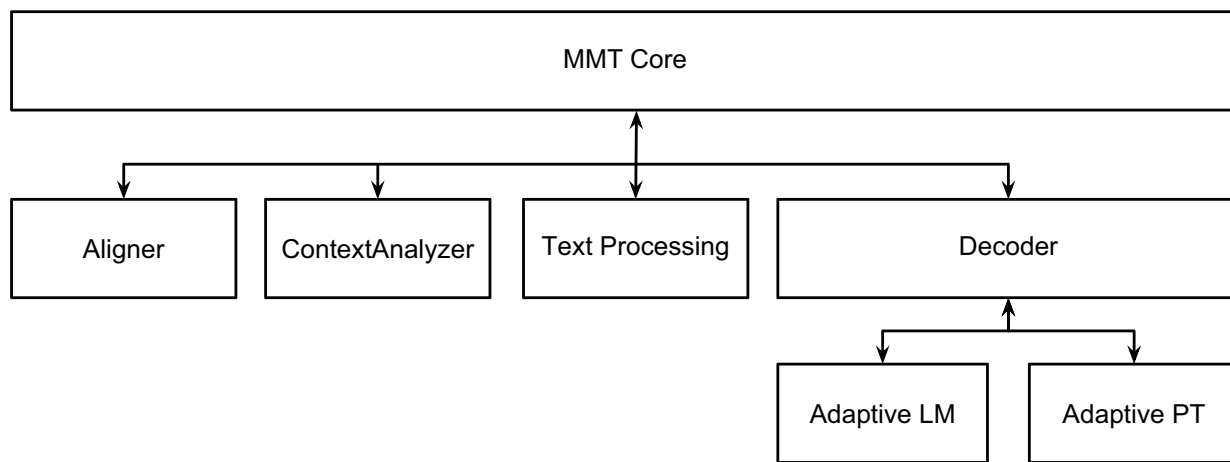


# Modern MT core components

---

Modern**MT** Next Generation  
Machine Translation

# Overall Architecture



# Word Alignment

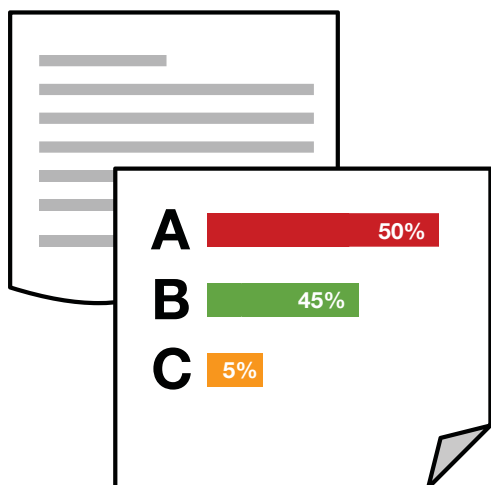
Object oriented **re-implementation** of FastAlign

Multithreading

Incremental training

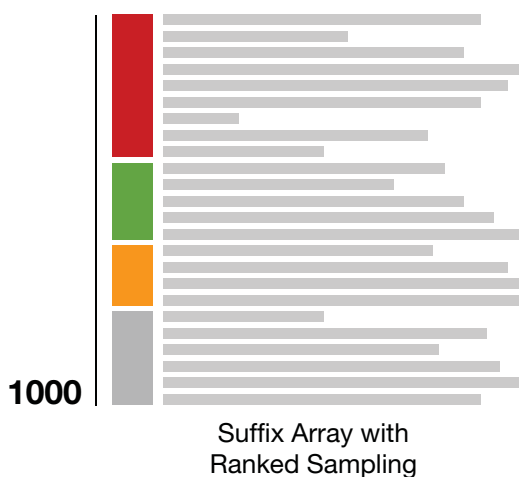
Giza++	FastAlign++
48,000 sec	2,800 sec ( <b>17x speed up</b> )
19.3 BLEU	18.9 BLEU ( <b>-0.4 loss</b> )

# Context Analyzer



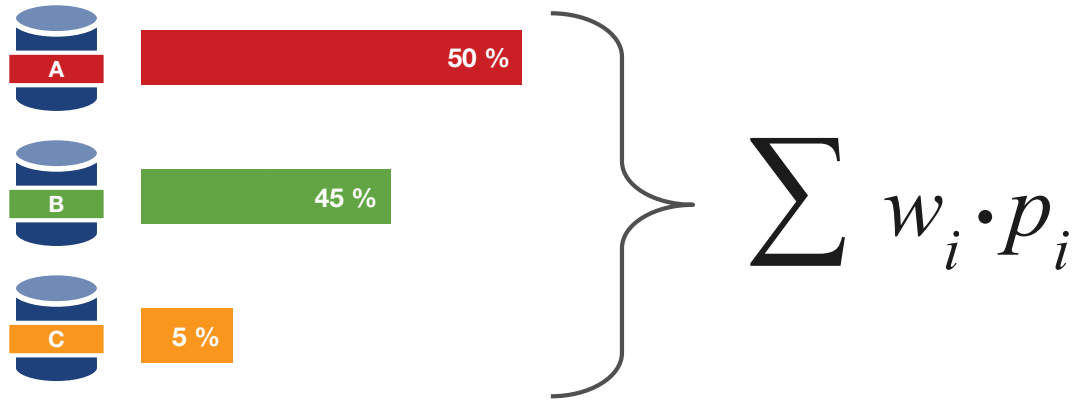
- Analyze the input text (tokenization, stop words)
- Retrieves best matching TMs
- Computes matching score

# Adaptive Phrase Table



- Suffix array indexed with TMs
- Phrase table is built on the fly by sampling from the SA
- Phrases of TMs with highest weights sampled first

# Adaptive Language Model



ModernMT Next Generation Machine Translation



## Modern MT cool features

ModernMT Next Generation Machine Translation

# Cool features

- **Fast** processing
- **Integrated** in the translation process
- **45 languages** tokenizer
- Automatic TMs **cleaning**
- Built-in **expression localizer**
- **Tag projector** for automatic placement of XML tags



# Word Tokenizer

**One interface** to 8 open-source tokenizers

including **re-implementation** of Moses tokenizer

	Moses Perl Tokenizer	MMT Tokenizer
Languages	21	45 <b>(+24)</b>
Speed*	17k w/s	340k w/s <b>(x20)</b>

\* 4 CPU, 83M word English corpus

# TM Cleaning

- **Draft translations cleaning**

EN This is **just** an example → IT Questo è un esempio  
 IT Questo è **solo** un esempio

- **XML tags cleaning and text normalization**

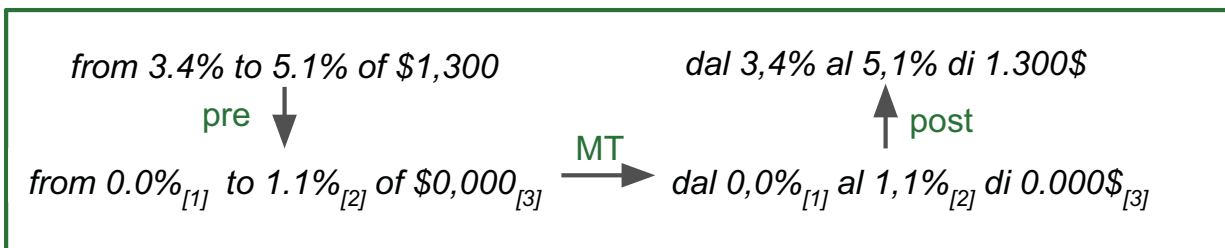
Did you know that the **<b>**Lunar distance**</b>** is **>** 300.000km ?  
 ↓  
 Did you know that the Lunar distance is **&gt;** 300.000km ?

# Numeric Expressions

Convert digits into placeholders

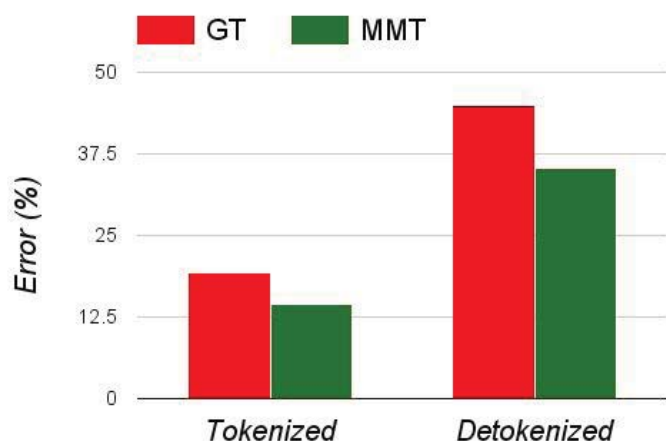
Translate with placeholders

Apply transformation and heuristics (for unaligned expr)





# Numeric Expressions



Test on 65 segments  
 135 numerical expressions  
 MMT better than GT  
 (rel. delta 25%-21%)

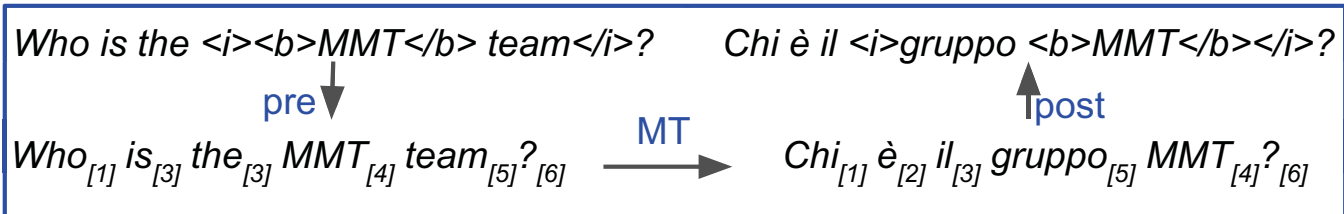
# Tag Manager

Identify, classify, and remove tags

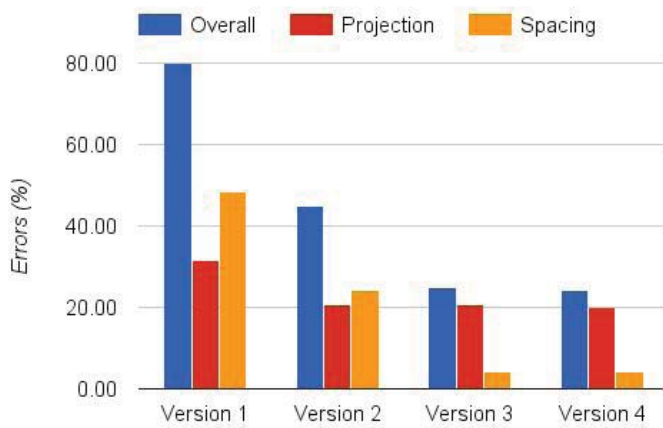
Translate w/o tags

Search insertion points using alignments and heuristics

Handle opening/closing, self-closing, nested, malformed tags



# Tag Manager



Tag projection error < 20%

Spacing errors around tags <= 4.2%



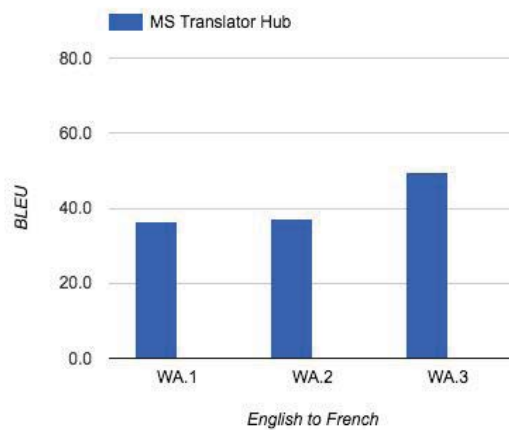
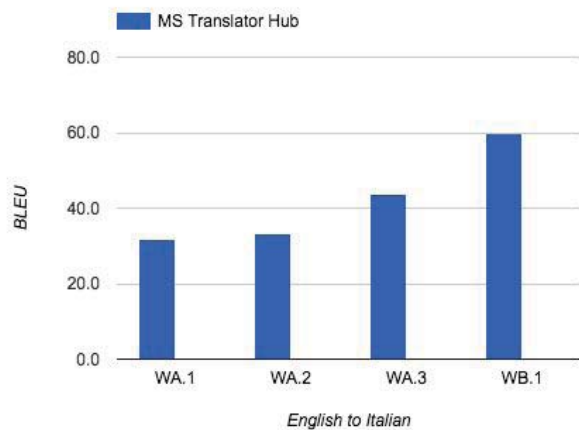
# Modern MT field testing

# Field testing

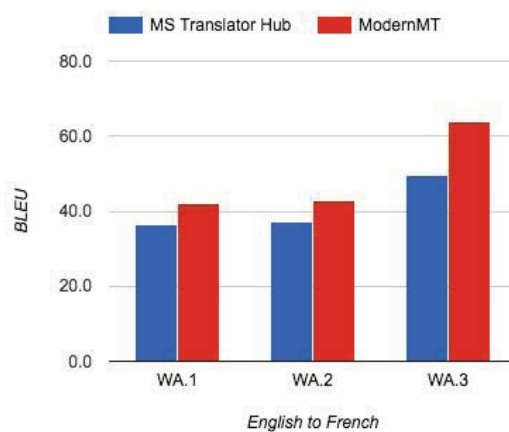
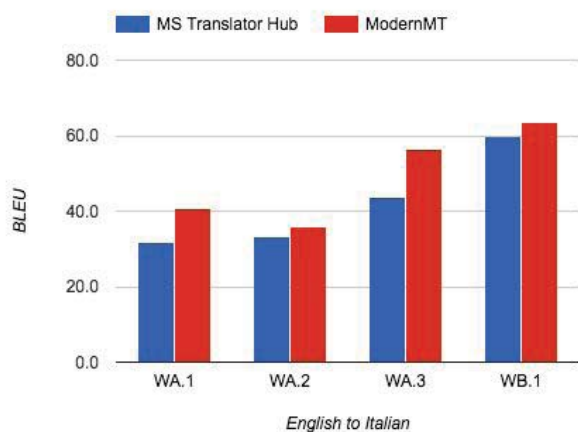
## Real business cases:

- Two world-class web companies
- Private TMs and real documents as test data
- Compared with current state-of-the-art **adaptive machine translation**.
- Two language pairs: **EN→IT** and **EN→FR**

# MS Translator Hub vs Modern MT



# MS Translator Hub vs Modern MT



**ModernMT** Next Generation Machine Translation

# Thank you

- easy to use product
- context aware translation
- incremental system
- open source, no licensing



[github.com/ModernMT/MMT](https://github.com/ModernMT/MMT)