# Language technologies for Education: recent results by the MLLP group

Alfons Juan
2nd Internet of Education Conference 2015
18 September 2015, Sarajevo

# Contents

- Research group at the *Univ. Politècnica de València* (Spain)

- **Research areas:**

  – Machine Learning and Applications

  – Natural Language Processing

  – Educational Technologies and Big Data

- **Recent research projects:**

  – trans**Lectures**: Transcription and translation of video lectures

  – EMMA: European Multiple MOOC Aggregator

  – Active Interaction for Speech Transcription and Translation

# transLectures (Nov 2011 – Oct 2014)

**Main goal:** to develop innovative, cost-effective (quasi-automatic) solutions to produce accurate subtitles for educational videos

**Two pilots:**

- *VideoLectures.NET:*   En, Sl   En→{De,Sl,Fr,Es}   Sl→En
- *poliMedia:*   Es, Ca   Es↔En   Ca↔{Es,En}

**Three scientific and technological objectives:**

- *Massive adaptation* to improve subtitling quality

- *Intelligent interaction* to improve subtitling quality

- *Integration into Opencast* to enable real-life evaluation

# trans**Lectures**: VideoLectures.NET



$>$20000 videos ($45$ min. on avg.): 85% English, 13% Slovenian, . . .

# trans**Lectures**: **poliMedia**



27000 videos (2-10 min.): 88% Spanish, 3% Catalan, . . .

# transLectures: Massive adaptation



- Key advances: neural networks and "extra" resources (slides)
- ASR: all error rates below 30 % (Sl: 26.9 %, others below 20 %)
- SMT: all BLEU scores above 20 (En→Sl: 20.1, mostly above 25)

# transLectures: Intelligent interaction



***poliMedia transcription***

The 10% of words recognized with minimum confidence include 30% of the actual word recognition errors

# transLectures: Integration into Opencast



>300 different institutions worldwide

# transLectures: try our tools



Try our tools at  `ttp.mllp.upv.es`

190 users and 1197 videos (264 hours) since May 2014

# transLectures: try our tools (cont.)

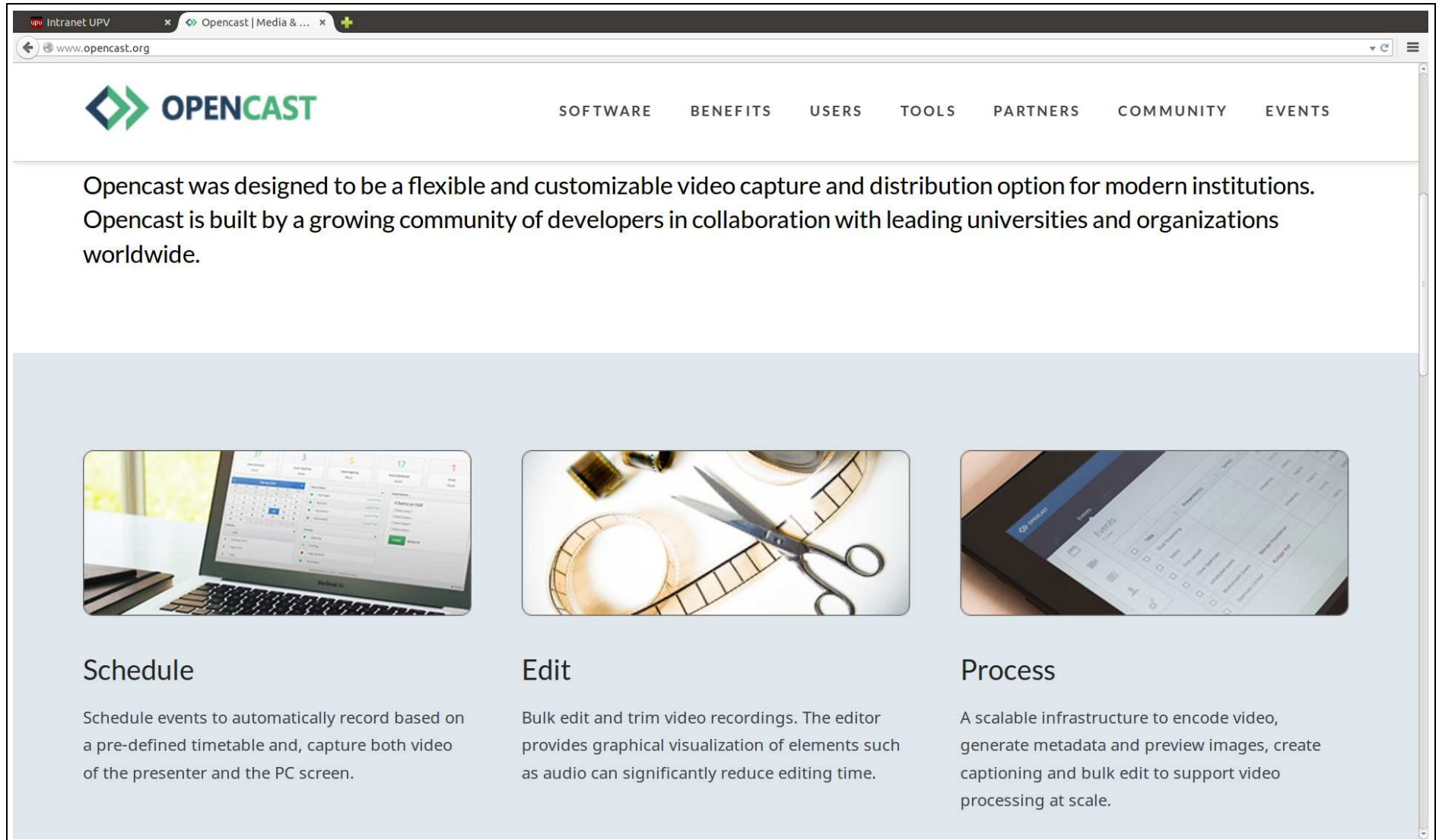| Institution | Country | Language |
|---|---|---|
| Univ. Politècnica de València | Spain | ES, CA, EN |
| Univ. Carlos III de Madrid | Spain | ES, EN |
| Universidade Aberta | Portugal | PT, EN |
| University of Naples Federico II | Italy | IT, EN |
| Open Universiteit Nederland | Netherlands | NL, EN |
| University of Leicester | UK | EN, ES |

**Examples:**

- Lecture recorded at poliMedia studios. [link]

- Lee Rubenstein, edX VP of Business Development: "Reinventing education", 30 June 2015, València. [link]

- Text-to-Speech demo [link]

# EMMA (Feb 2014 – Jul 2016)

## *European Multiple MOOC Aggregator (EMMA)*

**Main goal:** to provide multilingual access to European MOOCs

**Motivation:**

- Most MOOCs are offered in few languages (En, Es, Fr)

- Language barrier is keeping many learners from taking MOOCs

- MOOC components: texts, images, videos, forums

- EMMA uses trans**Lectures** tools to translate videos and texts
  - Few hours of video in 7 En, Es, It, Nl, Et, Pt and Fr
  - Source language is the national language of the MOOC provider
  - Target languages: En, Es and It

# EMMA: cost of manually translating MOOCs

**Texts:**

- Manual translation rate is approximately 2500 words per day
- A 6-week course with 75000 words takes 1.5 PM

**Videos:**

- Before translating, videos are manually transcribed (10 RTF)
- Then, transcriptions are translated (30 RTF)
- A course including 2 hours of video takes 0.5 PM

**Solutions to lower costs:**

- Crowdsourcing (e.g. TED talks)
- ASR and MT: user effort is reduced to 30% ($2 \rightarrow 0.6$ PM)

# EMMA: automatic video subtitling

1. Generation of *automatic transcriptions* from video

2. *Manual review* of automatic transcriptions



3. Generation of *automatic translations* from transcriptions

4. *Manual review* of automatic translations

# EMMA: automatic document translation

- Course text is ingested into the translation system

- Source and target texts are reviewed in parallel

- Preview of source and target texts also available

- Translated text is imported back into the EMMA platform

# EMMA: evaluations

## Video subtitling

| Language pairs | Transcription RTF (10) | Translation RTF (30) | Total RTF (40) |
|---|---|---|---|
| Spanish $\rightarrow$ English | 3 | 7 | 10 |
| English $\rightarrow$ Spanish | 6 | 17 | 23 |

## Document translation

| Language pairs | Translation RTF |
|---|---|
| Spanish $\rightarrow$ English | 7 |
| English $\rightarrow$ Spanish | 17 |

# EMMA: conclusions

- Multilingual access to your course boosts visibility

- The cost of manually translating your course is high (2 PM)

- Automatic translation can reduce the temporal cost to 30%

- Accuracy of automatic translation depends on several factors:

  – Languages involved
  – Availability of annotated data resources related to your course
  – Specificity of the course

- Designing a multilingual MOOC should also take into account: slides, images, application interfaces (demos), bibliography