# Global Under-Resourced Media Translation
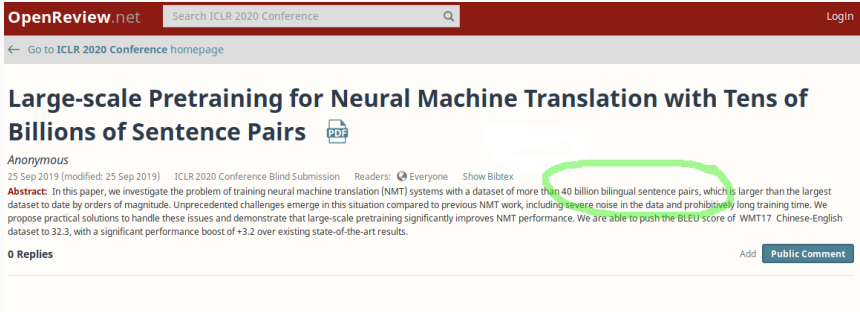
MetaForum

Barry Haddow

October 8th, 2019

# Sometimes You Have Lots of Data ...



But what if you don't have 40 billion parallel sentences?

## Global Under-Resourced Media Translation

**Rationale**

- MT is still poor for most world languages

**Aims**

- Improve translation quality
- Apply to journalism and media analyst use-cases



https://gourmet-project.eu/

## Techniques

- Data Gathering and Augmentation
- Modelling Morphological Structure
- Structure Induction at the Sentence Level
- Transfer Learning

## Research Partners



## User Partners

## Data Gathering and Augmentation

**Aims**

- Collection of corpora, lexical and linguistic resources (including web-scraping for low-resource)
- Data augmentation to extend corpora, applying rule-based techniques

**Progress**

- Delivered report on resources for all proposed languages
- Sponsored English-Gujarati task at WMT19
  - $\rightarrow$ Released new parallel and monolingual training corpoa
  - $\rightarrow$ Also development and test sets
- Sponsoring English-Tamil task at WMT20

## Modelling Challenges

**Challenges**

- Lack of data
- Structurally distant languages
- Morphological complexity and agglutination

**Approaches**

- Unsupervised and semi-supervised approaches
- Using linguistic and lexical resources
- Joint modelling of alignment and morphology
- Modelling of latent structure
- Better explanation of data via joint source-target modelling

## Transfer Learning: Low-resource Systems at WMT19

### English↔Gujarati

- Exploit large hi-en corpus
- Build unsupervised hi-gu system
- Using transliteration
  . . . and similarity
→ Synthesise new gu-en corpus

### English→Kazakh

- Exploit en-kk and kk-ru corpora
- Pivoted back-translation
- Continued training
- Hybridisation with RBMT

All systems scored highly in the official (human) evaluation

## GoURMET Translate

**Input**

> ગલ્ફ તણાવ : ઈરાનને નહીં રોકવામાં આવે તો દુનિયામાં તેલનો ભાવ વધશે - સાઉદી પ્રિન્સ સલમાન

Gujarati ▾

[Translate] [Clear]

**Output**

> Gulf tensions: If Iran is not stopped, the price of oil in the world will rise - Saudi Prince Salman

English ▾

- Models dockerised with secure API – enables integration of user tools
- First delivery (June 2019):
  - English ↔ Bulgarian, Gujarati, Swahili and Turkish
- Next delivery (March 2020)
  - English ↔ Amharic, Kyrgyz, Serbian and Tamil