



**Engaging Content**  
Engaging People

# Maintaining sentiment polarity in translation of user- generated content

Pintu Lohar, Haithem Afli and Andy Way

ADAPT Centre, School of Computing, Dublin City  
University

- Objective & Motivation
- Sentiment analysis of user-generated content
- Data Preparation
  - Corpus development
  - Sentiment annotation and classification
- Experiments
  - Sentiment Translation Architecture
  - Results
  - Discussion
- Conclusions and future work

- Analyse sentiment preservation & MT quality in the context of user-generated content (UGC)

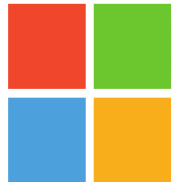
- Analyse sentiment preservation & MT quality in the context of user-generated content (UGC)
- Focus on whether sentiment classification helps improve sentiment preservation in MT of UGC

- Translation quality *per se* is not the main concern

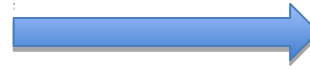
- Translation quality *per se* is not the main concern
  - **Sentiment preservation** is (arguably) more important

e.g. companies want to know what their customers think of their products and services.

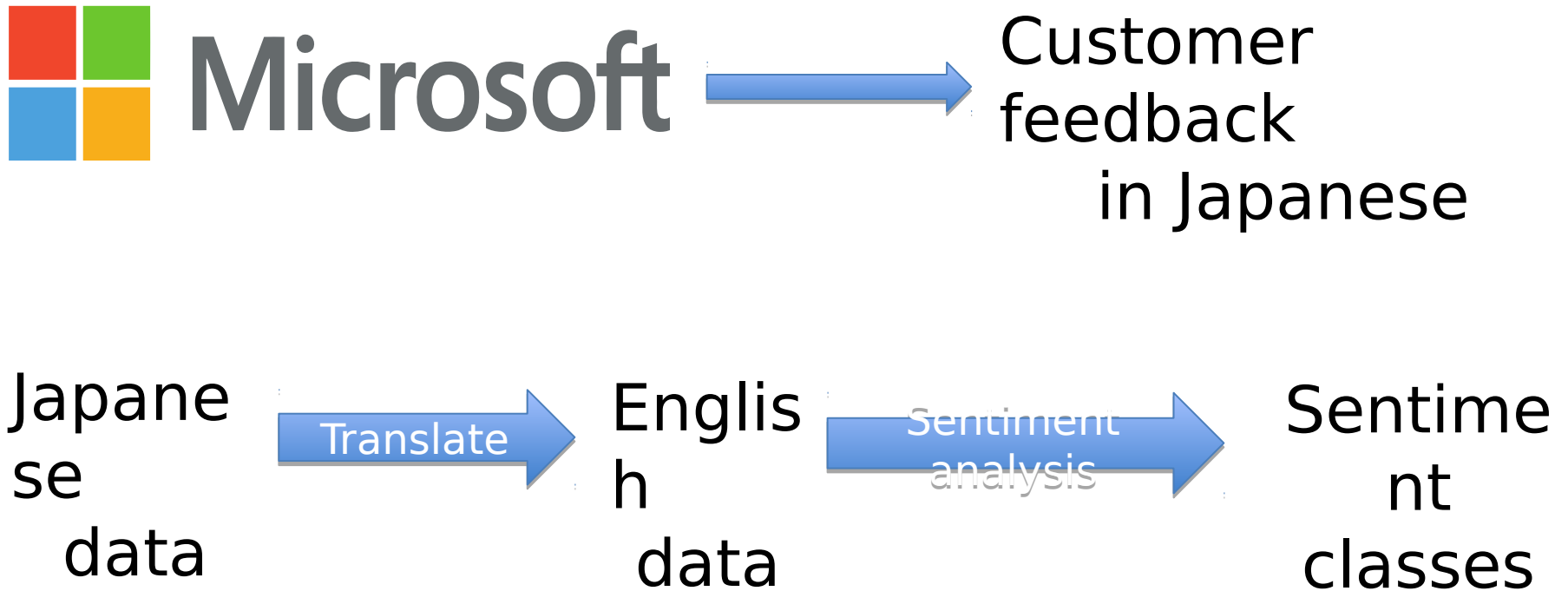
It is **crucial** that user sentiment in one language is preserved in the target language (typically, English).



Microsoft



Customer  
feedback  
in Japanese





# Track Record in UGC

www.adaptcentre.ie

## Filter Results

Number of Tweets: 835,725 / 835,725

Date ☐

Time ☐

Start Date:

1 Jan 2016

Start Time:

12:00 AM

End Date:

19 May 2017

End Time:

12:00 AM

Party ☐

AAA-PBP

Direct Democracy Ireland

Fianna Fáil

Fine Gael

Green Party

Independent

Independent Alliance

Irish Democratic Party

Labour

Renua

Sinn Féin

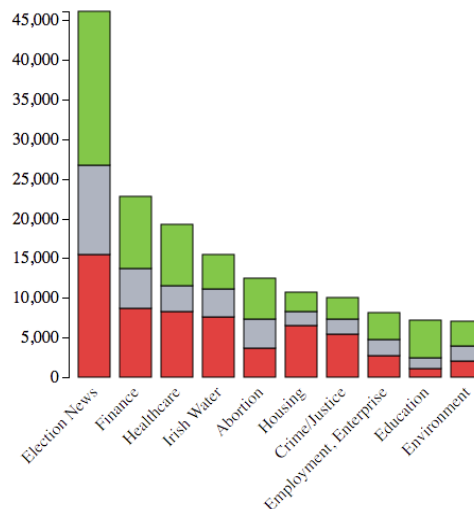


Home Timeline

Number of tweets: 835,725

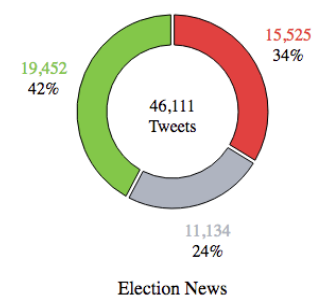
## Top Issues

Save



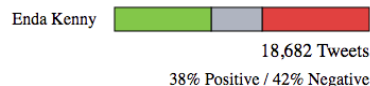
## Focused Issue

Save



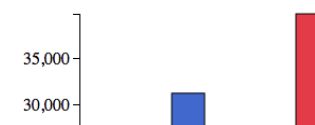
## Candidate Mentions

Save



## Top Party Mentions

Save



CNGL  
CENTRE FOR GLOBAL INTELLIGENT CONTENT

TRINITY  
COLLEGE  
DUBLIN

DCU

Microsoft

NDP  
National Development Plan

sfi  
Science Foundation  
Ireland



# Track Record in UGC

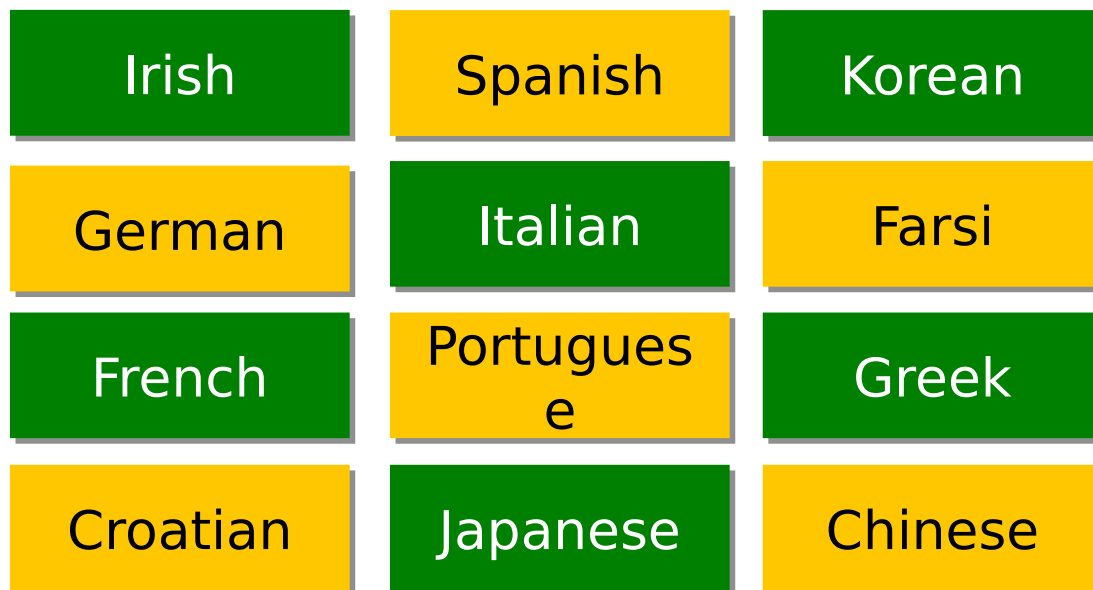
www.adaptcentre.ie

**13 languages and 24  
language pairs**



Brazilator

85,047,110 words in total



English

**CNGL**  
CENTRE FOR GLOBAL INTELLIGENT CONTENT



- UGC includes blog posts, podcasts, online videos, tweets etc.
- UGC is usually multilingual and of varying quality (sometimes deliberately)
- Sentiment analysis of UGC has many applications

# Related work


MT can alter the sentiment (Mohammad et al. (2016))



*Google Translate from English to German on 25/05/2017*

**English:** *he is out of the world cup*  
*negative*

**German:** *Er ist aus der weltmeisterschaft*  
*neutral*



- Can a sentiment classification approach help improve sentiment preservation in the target language ?

- Can a sentiment classification approach help improve sentiment preservation in the target language ?
- Is it useful to select a specific sentiment-MT model to translate the UGC with the same sentiment ?

## Corpus development:

- Twitter data set comprising 4,000 English tweets from the FIFA World Cup 2014 and their manual translations into German

## Corpus development:

- Twitter data set comprising 4,000 English tweets from the FIFA World Cup 2014 and their manual translations into German
- Informal translations of English tweets into German

e.g.                      English tweet  
German tweet

*Goaaaal*

*Toooooor*



## ▪ **Sentiment annotation**

Manually annotated sentiment scores  
between 0 and 1

## ▪ Sentiment annotation

Manually annotated sentiment scores between 0 and 1

## ▪ Sentiment classes

(i) **Negative**: sentiment score  $\leq 0.4$

(ii) **Neutral**: sentiment score  $\approx 0.5$

(iii) **Positive**: sentiment score  $\geq 0.6$

e.g.	Tweet	Sentiment score
------	-------	-----------------

	<i>injured Neymar out of World Cup</i>	0.2
--	--	-----

- Manual annotation of Twitter data is considered as the “gold-standard”

- Manual annotation of Twitter data is considered as the “gold-standard”
- 50 tweets per sentiment (negative, neutral and positive) are held out for tuning and testing purposes

Data	Train	Development			Test			Total
		#neg	#neu	#pos	#neg	#neu	#pos	
Twitter	3,700	50	50	50	50	50	50	4,000

*Data distribution of Twitter data for Training, development and test*

- Flickr and News commentary (``News'') data are used as additional resources
- Automatic sentiment analysis tool (Afli et. al. (2017)) is applied to Flickr and News data

- Flickr and News commentary (``News'') data are used as additional resources
- Automatic sentiment analysis tool (Afli et. al. (2017)) is applied to Flickr and News data

## **Performance accuracy:**

- 2,994 tweets out of 4,000 correctly classified by this tool when compared to the 'gold standard' data
- Accuracy = **74.85%**

Data	Sentiment classification	#neg	#neu	#pos	#total
Twitter	manual	919	1,308	1,473	3,700
Flickr	automatic	9,677	11,065	8,258	29,000
News	automatic	111,337	14,306	113,200	238,843

*Data distribution after sentiment classification*

## I. Translation without sentiment classification



## **I. Translation without sentiment classification**

## **II. Translation with sentiment classification**

- i. Manual sentiment classification (only Twitter data)
- ii. Automatic sentiment classification (Flickr & News data)

## **I. Translation without sentiment classification**

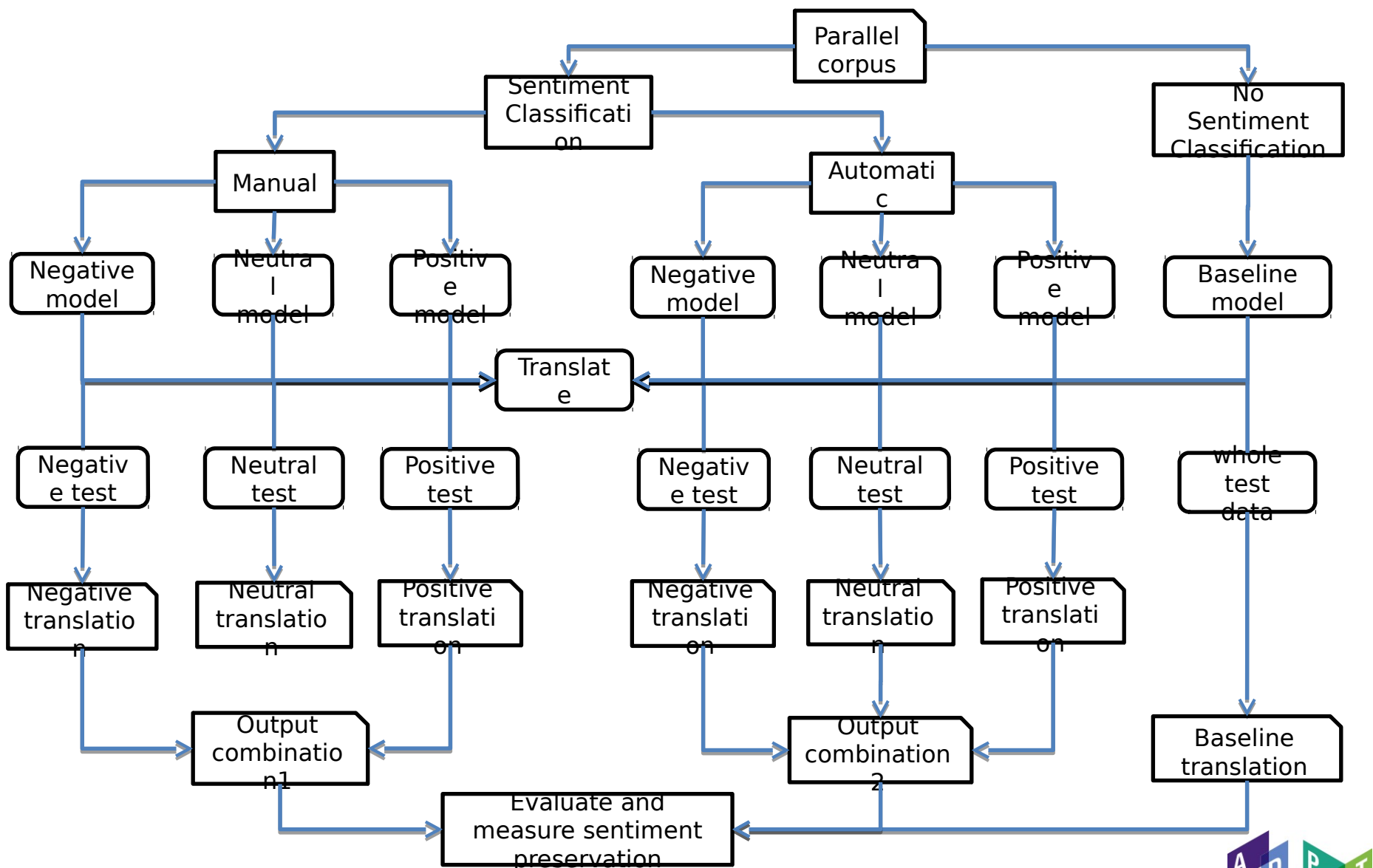
## **II. Translation with sentiment classification**

- i. Manual sentiment classification (only Twitter data)
- ii. Automatic sentiment classification (Flickr & News data)

## **III. Translation by wrong MT engines**

- i. Negative tweets by positive model
- ii. Neutral tweets by negative model
- iii. Positive tweets by neutral model

# Sentiment Translation Architecture



# Results

Translation model	Data size	Sentiment Classification	BLEU	METEOR	TER	Sentiment Preservation
Twitter	4k	✓	48.2	59.4	34.2	72.66%
Twitter (Baseline)		×	50.3	60.9	31.9	66.66%
Twitter + Flickr	33k	✓	48.5	59.8	33.9	71.33%
Twitter + Flickr		×	50.7	62.0	31.3	62.66%
Twitter + Flickr + News	272k	✓	50.3	62.3	31.0	<b>75.33%</b>
Twitter + Flickr + News		×	<b>52.0</b>	<b>63.4</b>	<b>30.1</b>	73.33%
Twitter (Wrong MT engine)	4k	✓	46.9	57.9	35.4	47.33%

*Experimental evaluation with data concatenation*



# Results

Translation model	Data size	Sentiment Classification	BLEU	METEOR	TER	Sentiment Preservation
Twitter	4k	✓	48.2	59.4	34.2	72.66%
Twitter (Baseline)		×	50.3	60.9	31.9	66.66%
Twitter + Flickr	33k	✓	48.5	59.8	33.9	71.33%
Twitter + Flickr		×	50.7	62.0	31.3	62.66%
Twitter + Flickr + News	272k	✓	50.3	62.3	31.0	<b>75.33%</b>
Twitter + Flickr + News		×	<b>52.0</b>	<b>63.4</b>	<b>30.1</b>	73.33%
Twitter (Wrong MT engine)	4k	✓	46.9	57.9	35.4	47.33%

*Experimental evaluation with data concatenation*



# Results

Translation model	Data size	Sentiment Classification	BLEU	METEOR	TER	Sentiment Preservation
Twitter	4k	✓	48.2	59.4	34.2	72.66%
Twitter (Baseline)		×	50.3	60.9	31.9	66.66%
Twitter + Flickr	33k	✓	48.5	59.8	33.9	71.33%
Twitter + Flickr		×	50.7	62.0	31.3	62.66%
Twitter + Flickr + News	272k	✓	50.3	62.3	31.0	<b>75.33%</b>
Twitter + Flickr + News		×	<b>52.0</b>	<b>63.4</b>	<b>30.1</b>	73.33%
Twitter (Wrong MT engine)	4k	✓	46.9	57.9	35.4	47.33%

*Experimental evaluation with data concatenation*



# Results

Translation model	Data size	Sentiment Classification	BLEU	METEOR	TER	Sentiment Preservation
Twitter	4k	✓	48.2	59.4	34.2	72.66%
Twitter (Baseline)		×	50.3	60.9	31.9	66.66%
Twitter + Flickr	33k	✓	48.5	59.8	33.9	71.33%
Twitter + Flickr		×	50.7	62.0	31.3	62.66%
Twitter + Flickr + News	272k	✓	50.3	62.3	31.0	<b>75.33%</b>
Twitter + Flickr + News		×	<b>52.0</b>	<b>63.4</b>	<b>30.1</b>	73.33%
Twitter (Wrong MT engine)	4k	✓	46.9	57.9	35.4	47.33%

*Experimental evaluation with data concatenation*



# Results

Translation model	Data size	Sentiment Classification	BLEU	METEOR	TER	Sentiment Preservation
Twitter	4k	✓	48.2	59.4	34.2	72.66%
Twitter (Baseline)		×	50.3	60.9	31.9	66.66%
Twitter + Flickr	33k	✓	48.5	59.8	33.9	71.33%
Twitter + Flickr		×	50.7	62.0	31.3	62.66%
Twitter + Flickr + News	272k	✓	50.3	62.3	31.0	<b>75.33%</b>
Twitter + Flickr + News		×	<b>52.0</b>	<b>63.4</b>	<b>30.1</b>	73.33%
Twitter (Wrong MT engine)	4k	✓	46.9	57.9	35.4	47.33%

*Experimental evaluation with data concatenation*





# Examples

Example	Reference	Sentiment translation model	Baseline model
1	Howard Web is a terrible ref #WorldCup	Howard Web is a schrecklicher ref #WorldCup	Howard Web is a schrecklicher ref #WorldCup
2	injured Neymar out of World Cup 2014	verletzter Neymar out the WC2014	verletzter Neymar out of World Cup 2014
3	penalty shootouts are too intense !	penalty shoot is to intensiv !	penalties is to intensiv !
4	damn chile is nice !!!! #WorldCup	freeking Chile is good !!! #WorldCup	damn Chile is good !!! #WorldCup
5	a bit boring ...	a little boring ...	some boring ...
6	im with Germany	I stand to Germany's side	I stand to Deutschlands side
7	as getting 1, GO CHILE !	completely matchless GO CHILE !	as getting 1, GO CHILE !



# Examples

Example	Reference	Sentiment translation system	Baseline system
1	Bosnia and Herzegovina really got f*** over man	Bosnia and Herzegovina eliminated echt demolished	Bosnia and Herzegovina were a abgezogen
2	when USA lost , but were still moving onto the next round	even if USA today we in the next round	could usa loses the next round
3	Brazil 5 WorldCup championship Argentina 2 WorldCup championship so Ill Examples where the sentiment is altered by the Baseline system	Brazil 5 time world champion Argentina 2 time world champion so Im for Brazil	Brazil 5 time world champions Argentina 2 time world champions so for Brazil

# Examples

Example	Reference	Sentiment translation system	Baseline system
1	Bosnia and Herzegovina really got f*** over man	Bosnia and Herzegovina eliminated echt demolished	Bosnia and Herzegovina were a abgezogen
2	when USA lost , but were still moving onto the next round	even if USA today we in the next round	could usa loses the next round
3	Brazil 5 WorldCup championship Argentina 2 WorldCup championship so Ill go with Brazil	Brazil 5 time world champion Argentina 2 time world champion so Im for Brazil	Brazil 5 time world champions Argentina 2 time world champions so for Brazil

*Examples where sentiment is altered by the Baseline system*

# Examples

Example	Reference	Right MT engine	Wrong MT engine
1	little break on the #WorldCup for an amazing #Wimbledon final!	small Pause from the #WorldCup for a amazing #Wimbledon final!	kleine Pause of the #WorldCup for a erstaunliches #Wimbledon final!
2	yes !!!!!	yes !!!!!	so !!!!!
3	a bit boring ...	a little boring ...	some was ...

*Comparison between sentiment polarities using the right and wrong MT engine*

- ❑ MT scores are better when no sentiment classification is used
- ❑ Sentiment classification approach performs better than for systems where it is switched off

# Discussion

Translation model	Sentiment Classification	BLEU	Sentiment Preservation
Twitter	✓	48.2	72.66% (+6%)
Twitter (Baseline)	×	50.3 (+2.1)	66.66%
Twitter + Flickr	✓	48.5	71.33% <b>(+8.67%)</b>
Twitter + Flickr	×	50.7 (+2.2)	62.66%
Twitter + Flickr + News	✓	50.3	<b>75.33%</b> (+2%)
Twitter + Flickr + News	×	<b>52.0</b> (+1.7)	73.33%

*MT quality VS sentiment preservation*

- ❑ In most cases, the Baseline system produces better outputs in terms of BLEU score
- ❑ In some cases, interestingly, sentiment classification approach produces better outputs

- ❑ In most cases, the Baseline system produces better outputs in terms of BLEU score
- ❑ In some cases, interestingly, sentiment classification approach produces better outputs
- ❑ Using specific sentiment-MT model to translate a text with the same sentiment is better in both ways

Translation model	Sentiment Classification	BLEU	Sentiment Preservation
Twitter (Right MT engine)	✓	48.2 (+1.3)	72.66% (+25%)
Twitter (Wrong MT engine)	✓	46.9	47.33%

*MT quality vs. sentiment preservation*



## **Sentiment translation system using nearest classes:**

1. Combine the negative- and neutral-sentimented parallel twitter data

## **Sentiment translation system using nearest classes:**

1. Combine the negative- and neutral-sentimented parallel twitter data
2. Combine the neutral- and positive-sentimented parallel twitter data

## Sentiment translation system using nearest classes:

1. Combine the negative- and neutral-sentimented parallel twitter data
2. Combine the neutral- and positive-sentimented parallel twitter data
3. Build the “negative-neutral” translation model using the data in **Step 1**

## Sentiment translation system using nearest classes:

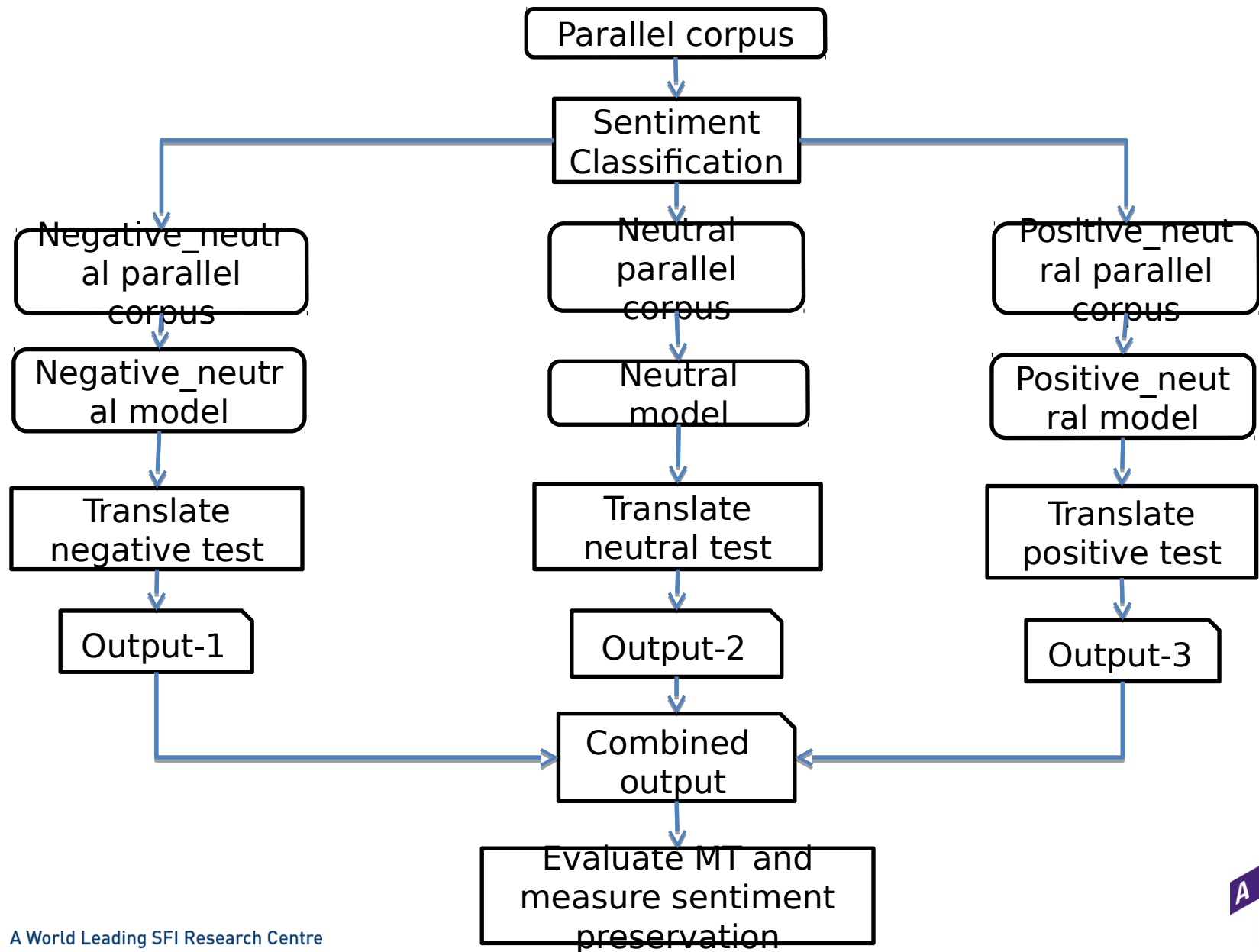
1. Combine the negative- and neutral-sentimented parallel twitter data
2. Combine the neutral- and positive-sentimented parallel twitter data
3. Build the “negative-neutral” translation model using the data in **Step 1**
4. Build the “positive-neutral” translation model using the data in **Step 2**

- Building a single translation system using the whole twitter data
- Building the sentiment-specific translation systems
- Building the nearest-sentiment translation systems

## Data distribution:

Experiments	Train	Development			Test			Total
		#neg	#neu	#pos	#neg	#neu	#pos	
Exp 1	3,700	50	50	50	50	50	50	4,000
Exp 2	3,400	100	100	100	100	100	100	4,000

# Architecture: nearest sentiment class



# Results

Exp1 (150 test data)	BLEU	Meteor	TER	Sentimen t Preservat ion
Twitter (Baseline)	<b>50.3</b>	<b>60.9</b>	<b>31.9</b>	66.66%
Twitter_SentCl ass	48.2 (- 2.1)	59.4 (- 1.5)	34.2 (+2.3)	<b>72.66%</b> (+6)
Twitter_NearS ent	49.0 (- 1.3)	60.1 (- 0.8)	34.0 (+2.1)	66.66% (+0)

# Results

Exp1 (150 test data)	BLEU	Meteor	TER	Sentimen t Preservat ion
Twitter (Baseline)	<b>50.3</b>	<b>60.9</b>	<b>31.9</b>	66.66%
Twitter_SentCl ass	48.2 (- 2.1)	59.4 (- 1.5)	34.2 (+2.3)	<b>72.66%</b> (+6)
Exp2 (300 test data)	BLEU	Meteor	TER	Sentimen t Preservat ion
Twitter (Baseline)	<b>51.3</b>	<b>62.5</b>	<b>31.0</b>	52.33%
Twitter_SentCl ass	47.3 (-4)	59.1 (- 3.4)	35.2 (+3.8)	<b>60.33%</b> (+8)
Twitter_NearS ent	48.3 (-3)	59.6 (- 2.9)	34.4 (+3.4)	60.0% (+7.67)





# Results

Exp1 (150 test data)	BLEU	Meteor	TER	Sentimen t Preservat ion
Twitter (Baseline)	<b>50.3</b>	<b>60.9</b>	<b>31.9</b>	66.66%
Twitter_SentCl ass	48.2 (- 2.1)	59.4 (- 1.5)	34.2 (+2.3)	<b>72.66%</b> (+6)
Exp2 (300 test data)	BLEU	Meteor	TER	Sentimen t Preservat ion
Twitter (Baseline)	<b>51.3</b>	<b>62.5</b>	<b>31.0</b>	52.33%
Twitter_SentCl ass	47.3 (-4)	59.1 (- 3.4)	35.2 (+3.8)	<b>60.33%</b> (+8)
Twitter_Nears ent	48.3 (-3)	59.6 (- 2.9)	34.4 (+3.4)	60.0% (+7.67)

Note that scores are different for Exp1 and Exp2 as they are tested on different data distributions



- ❑ A small amount of test data (i.e., 150 tweet pairs) is not (really) sufficient to confirm our hypothesis
- ❑ Comparatively larger test data (i.e., 300 tweet pairs) confirms the utility of our approach
- ❑ Our approach is able to reduce the gap between translation quality and sentiment preservation

- Despite a small deterioration in translation quality, our approach significantly improves sentiment preservation
- It is essential to carefully select the proper MT engine conveying the same sentiment polarity as that of the UGC
- The nearest sentiment class-inclusion approach helps improve the balance between MT quality and sentiment preservation

- Apply to other language pairs and also other forms of UGC such as customer feedback, blogs etc.
- Further refine the sentiment classes (strong positive, strong negative etc.) in order to build more specific translation models

# Thank you