

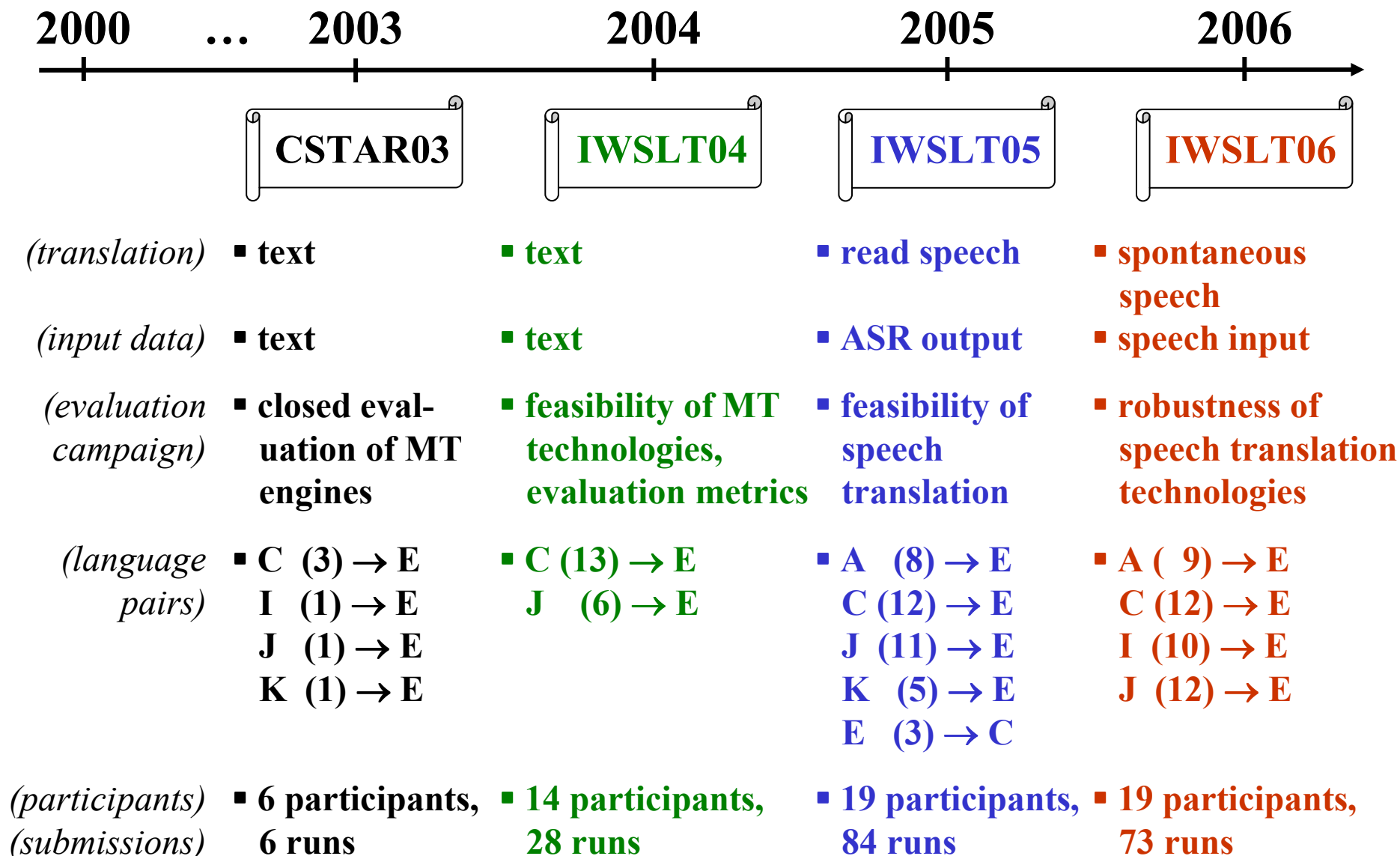
# Overview of the IWSLT 2006 Evaluation Campaign



Michael Paul

ATR Spoken Language Communication Research Laboratories  
Kyoto, Japan

# History of IWSLT



# Outline of Talk

## 1. Evaluation Campaign:

- data preparations
- translation input and data track conditions
- run submissions
- evaluation specifications

## 2. Evaluation Results:

- subjective/automatic evaluation
- correlation between evaluation metrics

## 3. Discussions:

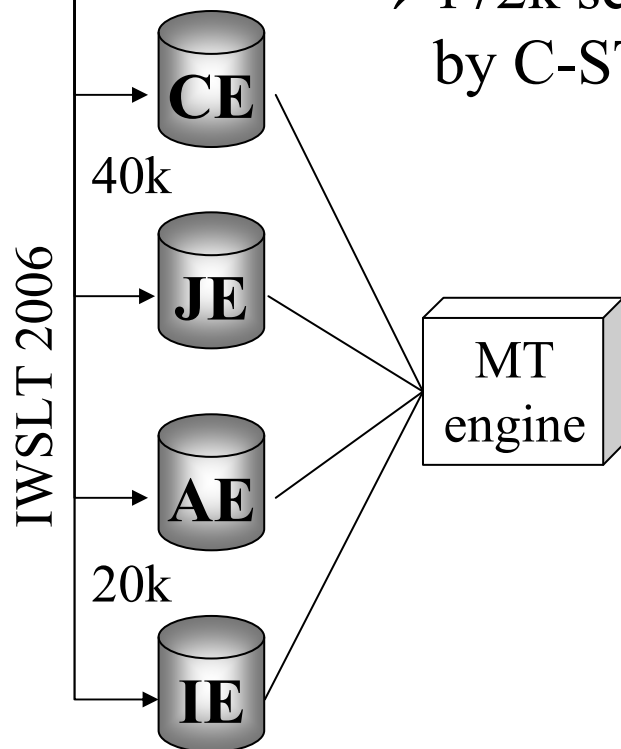
- Challenge Task 2006
- source language effects
- robustness towards recognition errors
- innovative idea's explored by participants

## BTEC



→ useful sentences, together with the translation into other languages usually found in phrasebooks for tourists going abroad

→ 172k sentence pairs collected/translated by C-STAR partners (A,C,E,I,J,K)



**J:** フィルムを買いたいです。

**E:** I want to buy a roll of film.

**J:** 8人分予約したいです。

**E:** I 'd like to reserve a table for eight.

**J:** 友人が車にひかれ大けがをしました。

**E:** My friend was hit by a car and badly injured.

# Supplied Resources

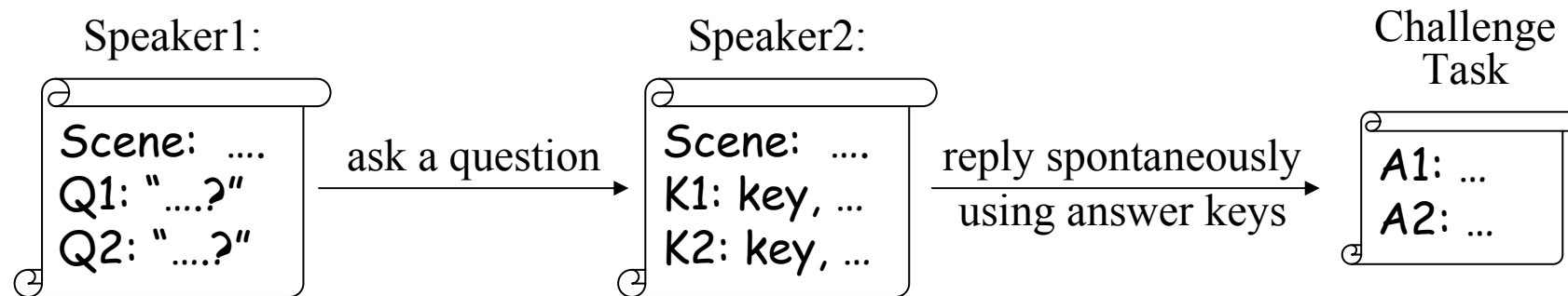


<i>type</i>	<i>language</i>	<i>sentence count</i>	<i>word token</i>	<i>word type</i>	<i>words per sentence</i>
training	<b>C/E</b>	40k	342k / 367k	11k / 7k	<b>8.6 / 9.2</b>
	<b>J/E</b>		398k / 367k	11k / 7k	<b>10.0 / 9.2</b>
	<b>A/E</b>	20k	154k / 183k	18k / 5k	<b>7.7 / 9.2</b>
	<b>I/E</b>		171k / 183k	10k / 5k	<b>8.6 / 9.2</b>
development (dev1, dev2, dev3)	<b>C/E<sub>16</sub></b>	1.5k	11k / 198k	2k / 3k	<b>7.0 / 8.2</b>
	<b>J/E<sub>16</sub></b>		12k / 198k	2k / 3k	<b>8.2 / 8.2</b>
	<b>A/E<sub>16</sub></b>		9k / 198k	3k / 3k	<b>6.3 / 8.2</b>
	<b>I/E<sub>16</sub></b>		10k / 198k	2k / 3k	<b>6.8 / 8.2</b>

# Challenge Task 2006



- speech input with certain level of “spontaneity”
- **spontaneous answers** to questions in tourism domain



[airplane] passenger asks flight attendance for help

Q1: Okay. Where can I put my luggage? Is it here okay?

K1: *[not here],  
[overhead compartment]*

A1: **sorry you'd better put it in the overhead compartment**

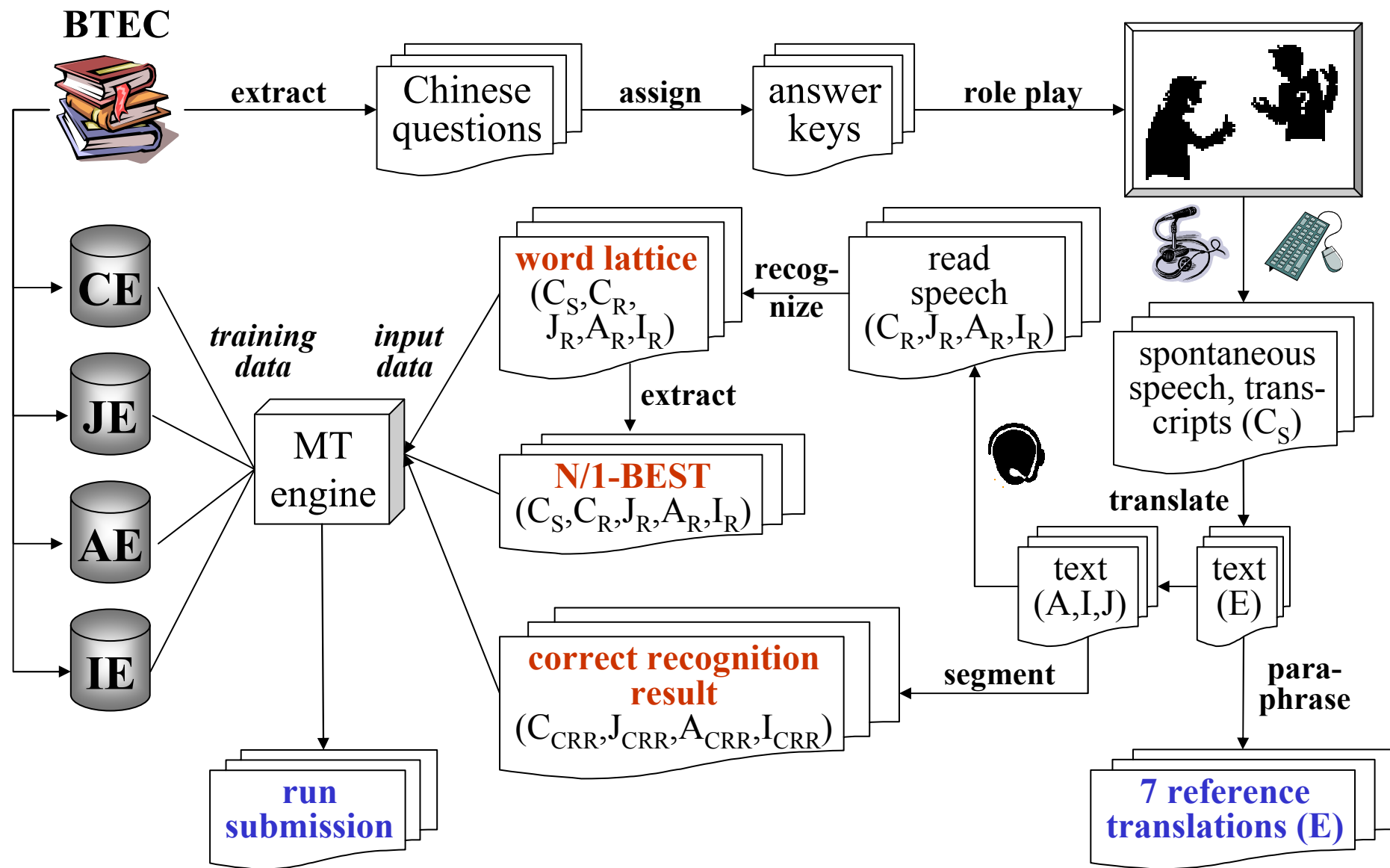
[airport] customer asks taxi driver for directions

Q2: Take me to this address. How long will it take?

K2: *[depending on traffic condition],  
[around 20 minutes]*

A2: **it's hard to say it depends on the traffic condition it should take only twenty minutes or so if there's no traffic jam**

# Data Preparation



# Data Statistics



<i>task</i>	<i>language</i>	<i>sentence count</i>	<i>word token</i>	<i>word type</i>	<i>words per sentence</i>
eval	<b>C/E<sub>7</sub></b>	500	6.0k / 50k	1.3k / 1.6k	<b>12.1 / 14.4</b>
	<b>J/E<sub>7</sub></b>		7.4k / 50k	1.2k / 1.6k	<b>14.8 / 14.4</b>
	<b>A/E<sub>7</sub></b>		5.2k / 50k	1.9k / 1.6k	<b>10.4 / 14.4</b>
	<b>I/E<sub>7</sub></b>		6.7k / 50k	1.5k / 1.6k	<b>13.4 / 14.4</b>

<i>task</i>	<i>Out-Of-Vocabulary rates (%)</i>			
	<i>CRR</i>	<i>1BEST</i>	<i>NBEST</i>	
eval	<b>C<sub>S</sub></b>	2.6	2.1	2.4
	<b>C<sub>R</sub></b>	2.6	2.4	2.5
	<b>J<sub>R</sub></b>	2.2	1.6	2.3
	<b>A<sub>R</sub></b>	<b>14.3</b>	<b>16.0</b>	<b>17.1</b>
	<b>I<sub>R</sub></b>	4.3	2.5	2.6
	<b>E<sub>7</sub></b>	2.7 (for AE,IE) / 1.9 (for CE/JE)		



# Recognition Accuracy



<i>task</i>	<i>word (%)</i>		<i>sentence (%)</i>		
	<i>lattice</i>	<i>1BEST</i>	<i>lattice</i>	<i>1BEST</i>	
<b>eval</b>	<b>C<sub>S</sub></b>	79.08	<b>68.11</b>	22.80	<b>16.60</b>
	<b>C<sub>R</sub></b>	82.07	<b>73.64</b>	28.40	<b>22.80</b>
	<b>J<sub>R</sub></b>	90.48	<b>85.14</b>	52.60	<b>38.00</b>
	<b>A<sub>R</sub></b>	88.20	<b>73.88</b>	41.60	<b>16.60</b>
	<b>I<sub>R</sub></b>	72.90	<b>70.88</b>	5.40	<b>4.60</b>

- performance differences between source languages
  - closed language model used for Arabic and Chinese
- differences between lattice and 1BEST accuracies

## Translation Directions

- Arabic
  - Chinese
  - Italian
  - Japanese
- } → English

## Data Tracks

- **OPEN:**
  - in-domain training data restricted to supplied BTEC resources
- **CSTAR:**
  - no restrictions

## Input Conditions

### **Cleaned Transcripts**

*plain text*

(text normalized according to ASR engine)



correct recognition results of supplied ASR engines

### **ASR Output**

*word lattice*

*NBEST*

*1BEST*



output of ASR engine supplied by CSTAR

partner

### ~~**Speech Input**~~

~~*audio data*~~

~~(C: spontaneous speech)~~

~~(A,C,I,J: read speech)~~



~~each participant uses its own ASR engine~~

# Participants



Research Group		System	Type	Input
US	AT&T Research	ATT	SMT	$C_S, C_R$
EU	CLIPS-GETA	CLIPS	RBMT	$C_S^*, J_R^*, A_R^*, I_R^*$
EU	Dublin City University	DCU	EBMT	$A_R, I_R^*$
ZH	Hong Kong University	<b>HKUST</b>	SMT	<b><math>C_S, C_R, J_R, A_R, I_R</math></b>
US	IBM	IBM	SMT	$A_R$
EU	Instituto Trentino di Cultura	<b>ITC-irst</b>	SMT	<b><math>C_S, C_R, J_R, A_R, I_R</math></b>
US/EU	JHU.Summer Workshop 2006	<b>JHU_WS06</b>	SMT	$C_S, C_R$
JP	Kyoto University	<b>Kyoto-U</b>	EBMT	$J_R$
US	MIT Lincoln Lab / Air Force Research Lab	<b>MIT-LL-AFRL</b>	SMT	$C_S, C_R, J_R, I_R$
JP	National Institute of Science & Technology	<b>NAIST</b>	SMT	$J_R$
JP	NICT / ATR-SLC	<b>NiCT-ATR</b>	SMT	<b><math>C_S, C_R, J_R, A_R, I_R</math></b>
ZH	NLPR, Chinese Academy of Science	<b>NLPR</b>	RBMT, SMT	$C_S, C_R$
JP	NTT Communication Research	<b>NTT</b>	SMT	<b><math>C_S, C_R, J_R, A_R, I_R</math></b>
EU	Rheinisch Westfälische Hochschule	<b>RWTH</b>	SMT	$C_S, C_R, J_R$
EU	SHARP Laboratories of Europe	<b>SLE</b>	EBMT	$J_R$
EU	TALP-UPC Research Center (2x)	<b>TALP</b>	SMT	$C_R, J_R, A_R, I_R$
US	InterACT, CMU / Karlsruhe University (2x)	<b>UKACMU</b>	SMT	<b><math>C_S, C_R, J_R, A_R, I_R</math></b>
US	University of Washington	<b>Washington-U</b>	SMT	$I_R$
ZH	Xiamen University	<b>Xiamen-U</b>	SMT	$C_S, C_R$

# Run Submissions



Type	Input	Lang	Group	Runs OPEN / C-STAR
text	correct recognition result	$A_{CRR}$ $C_{CRR}$ $I_{CRR}$ $J_{CRR}$	19	<i>mandatory for all participants</i>
read speech	ASR output	$A_R$ $C_R$ $I_R$ $J_R$	9 12 10 12	11 (14) / 1 (1) 14 (17) / 3 (3) 12 (14) / 1 (3) 14 (14) / 2 (3)
spontaneous speech	ASR output	$C_S$	12	12 (11) / 3 (3)
<b>TOTAL</b>			<b>19</b>	<b>63 (70) / 10 (13)</b>

# Evaluation Specifications



- **case-sensitive, with punctuation marks** (*official*)
- case-insensitive, without punctuation marks (*additional*)

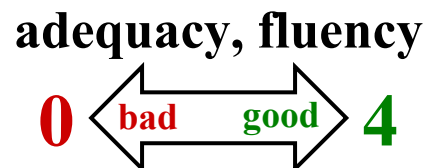
## Automatic Evaluation:

- all run submissions
- metrics:



## Subjective Evaluation:

- CE, ASR Output, Open Data Track
- 7 MT engines with  $C_S, C_R, C_{CRR}$  run submissions
- median of 3 human grades
- metrics:



### adequacy

4	All Information
3	Most Information
2	Much Information
1	Little Information
0	None

### fluency

4	Flawless English
3	Good English
2	Non-native English
1	Disfluent English
0	Incomprehensible

# Outline of Talk

## 1. Evaluation Campaign:

- data preparations
- translation input and data track conditions
- run submissions
- evaluation specifications

## 2. Evaluation Results:

- **subjective/automatic evaluation**
- **correlation between evaluation metrics**

## 3. Discussions:

- Challenge Task 2006
- source language effects
- robustness towards recognition errors
- innovative idea's explored by participants



## Subjective Evaluation

- *adequacy/fluency* : p.11 (scores, system rankings)

## Automatic Evaluation

- *BLEU/NIST/METEOR*: pp.12-13 (scores), pp.14-15 (system rankings)
  - test significance of differences in translation quality between two MT systems using “**bootStrap**” method:
    - (1) perform a random sampling with replacement from the *eval* data
    - (2) calculate respective evaluation metric scores of each MT engine and differences between the two MT engine scores
    - (3) repeat sampling/scoring steps iteratively
    - (4) apply Student’s t-test at a significant level of 95% to test whether score differences are significant
- **horizontal lines omitted in ranking tables, if system performance difference is NOT significant**

## TOP Scores (MT Engine)

metric	$C_{CRR}$	$C_R$	$C_S$
adequacy	<b>1.4319</b> (MIT-LL-AFRL)	<b>1.0297</b> (JHU-WS06)	<b>0.9647</b> (JHU-WS06)
fluency	<b>1.6498</b> (RWTH)	<b>1.2952</b> (RWTH)	<b>1.3734</b> (MIT-LL-AFRL)

## Combination of Subjective Evaluation Rankings

$C_{CRR}$	$C_R$	$C_S$
<b>MIT-LL-AFRL</b>	<b>JHU-WS06</b>	<b>JHU-WS06</b>
RWTH	<b>MIT-LL-AFRL</b>	RWTH
NTT	RWTH	NTT
<b>JHU-WS06</b>	NTT	<b>MIT-LL-AFRL</b>
NiCT-ATR	NiCT-ATR	UKACMU_SMT
UKACMU_SMT	UKACMU_SMT	NiCT-ATR



# Automatic Evaluation Results



## TOP Scores (MT Engine) for ASR Output

input	BLEU	NIST	METEOR
$C_S$	<b>0.1898</b> (RWTH)	<b>5.1513</b> (JHU-WS)	<b>0.4238</b> (HKUST)
$C_R$	<b>0.2111</b> (RWTH)	<b>5.4154</b> (MIT-LL-AFRL)	<b>0.4456</b> (HKUST)
$J_R$	<b>0.2142</b> (RWTH)	<b>5.6502</b> (RWTH)	<b>0.4574</b> (NiCT-ATR)
$A_R$	<b>0.2274</b> (IBM)	<b>5.9216</b> (NiCT-ATR)	<b>0.4867</b> (NiCT-ATR)
$I_R$	<b>0.2989</b> (NiCT-ATR)	<b>6.9318</b> (Washington-U)	<b>0.5853</b> (Washington-U)

# Automatic Evaluation Results



## Combination of Automatic Evaluation Rankings

$C_S$	$C_R$	$J_R$	$A_R$	$I_R$
RWTH	RWTH	RWTH	IBM	Washington-U
JHU-WS06	MIT-AFRL	NICT-ATR	NICT-ATR	NICT-ATR
NICT-ATR	NICT-ATR	UKACMU	TALP-tuples	TALP-tuples
UKACMU	JHU-WS06	NTT	TALP-comb	MIT-AFRL
HKUST	ITC-irst	MIT-AFRL	NTT	TALP-comb
ITC-irst	TALP-tuples	ITC-irst	UKACMU	ITC-irst
MIT-AFRL	TALP-phrases	SLE	TALP-phrases	TALP-phrases
NTT	UKACMU	HKUST	ITC-irst	NTT
Xiamen-U	HKUST	TALP-tuples	DCU	DCU
ATT	TALP-comb	NAIST	HKUST	UKACMU
NLPR	NTT	Kyoto-U	CLIPS	HKUST
CLIPS	Xiamen-U	TALP-comb		CLIPS
	NLPR	TALP-phrases		
	ATT	CLIPS		

# Correlation between Automatic and Subjective Evaluation Metrics

input	metric	BLEU	NIST	METEOR
$C_{CRR}$	fluency	<b>0.96</b>	0.84	0.93
	adequacy	0.95	0.82	<b>0.96</b>
$C_R$	fluency	<b>0.89</b>	0.63	0.66
	adequacy	0.83	0.64	<b>0.89</b>
$C_S$	fluency	<b>0.88</b>	0.55	0.72
	adequacy	0.34	<b>0.57</b>	0.54

# Outline of Talk

## 1. Evaluation Campaign:

- data preparations
- translation input and data track conditions
- run submissions
- evaluation specifications

## 2. Evaluation Results:

- subjective/automatic evaluation
- correlation between evaluation metrics

## 3. Discussions:

- **Challenge Task 2006**
- **source language effects**
- **robustness towards recognition errors**
- **innovative ideas explored by participants**

# Challenge Task 2006



- more difficult than previous IWSLT tasks

<i>translation task</i>	<i>training data</i>	
	<i>40k (CE/JE)</i>	<i>20k (AE/IE)</i>
<b>dev<sub>1</sub> / dev<sub>2</sub> / dev<sub>3</sub></b>	27.5 / 31.4 / 32.9	32.6 / 36.7 / 38.8
<b>dev<sub>4</sub> / eval</b>	<b>85.6 / 105.9</b>	<b>98.3 / 113.9</b>

- quite low MT performance for all systems for all conditions
  - discrepancy between training and evaluation data
  - high OOV figures
  - number of reference translations differed (16 vs. 7)

# Source Language Effects



## Italian:

- highest scores despite worst recognition accuracy  
→ close language relationship

## Arabic:

- largest OOV rates  
→ re-segmentation led to improved coverage & translation quality

## Japanese:

- highest recognition accuracy, but low scores  
→ one of the most difficult translation tasks
- largest number of non-SMT run submissions

## Chinese:

- recognition accuracy similar to Arabic, but much lower scores
- largest number of participants

**task complexity:**

**CE ≈ JE > AE ≫ IE**

# Robustness Towards Recognition Errors

TOP-scoring systems		recognition errors			
		none	low	medium	high
$C_S$	%	16.6	39.6	34.4	9.4
	adequacy	1.52	1.14	0.69	0.27
	fluency	1.81	1.22	1.21	0.58
$C_R$	%	22.8	45.5	22.4	9.3
	adequacy	1.80	1.02	0.63	0.17
	fluency	2.05	1.13	0.91	0.30

# Innovative Ideas Explored by Participants

- **additional training resources**
  - *in-domain* → large gain (CSTAR data track)
  - *out-of-domain* → partially effective (IBM, Washington-U)
- **distortion modeling** (ITC-irst, TALP)
- **topic-dependent model adaptation** (NiCT-ATR)
- **efficient decoding of word lattices** (JHU\_WS06, ITC-irst)
- **rescoring/ranking features** (NTT, RWTH, Washington-U)

**closer coupling of ASR and MT technologies required  
to overcome problems of speech translation tasks**



# The End

谢谢大家注意  
听我的发言。

ご静聴、ありがとう  
ございました。

**Thank you  
for your attention!**

شكراً على انتباهكم.

**Grazie  
per la vostra  
attenzione!**