

## Systems

Alina Karakanta, Francois Buet, Mauro Cettolo, Francois Yvon

HUMANE AI NET

### Problem

Reference-based evaluation for **perfect/imperfect** texts

**Per:** He said, <eob> "It almost <eob> doesn't matter what you know."

**Ref:** He said, <eob> "It almost doesn't matter <eol> what you know."

**Imp:** He told me, <eob> "What you know <eol> is not important."

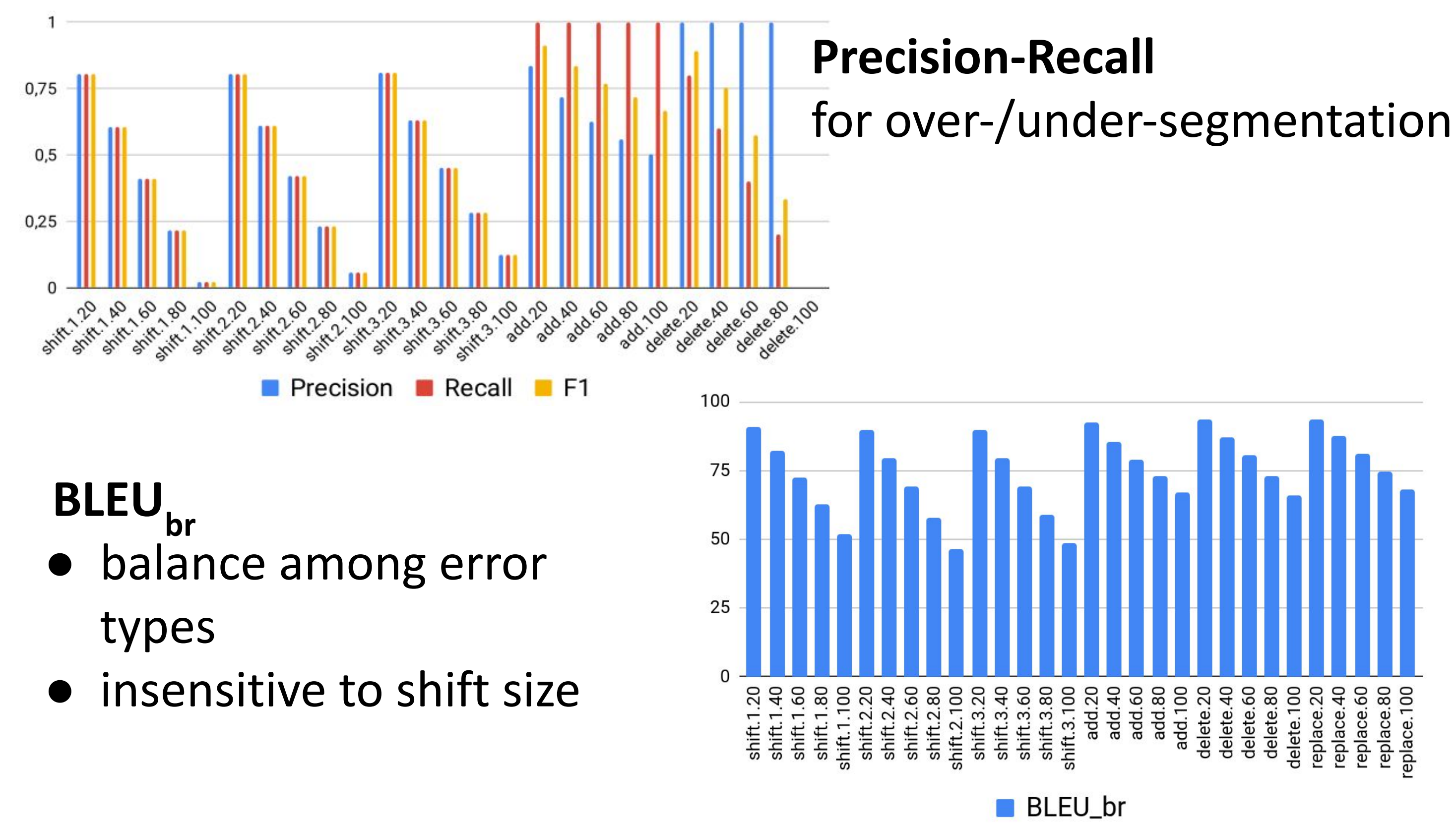
**Perfect: Precision=1/2 Recall=1/2 Imperfect: Precision? Recall?**

➔ Standard segmentation metrics cannot be computed for **Imperf**

**How to evaluate segmentation for Imperfect texts?**

### Exp1: Metric sensitivity/robustness

=> Controlled segmentation degradations of reference



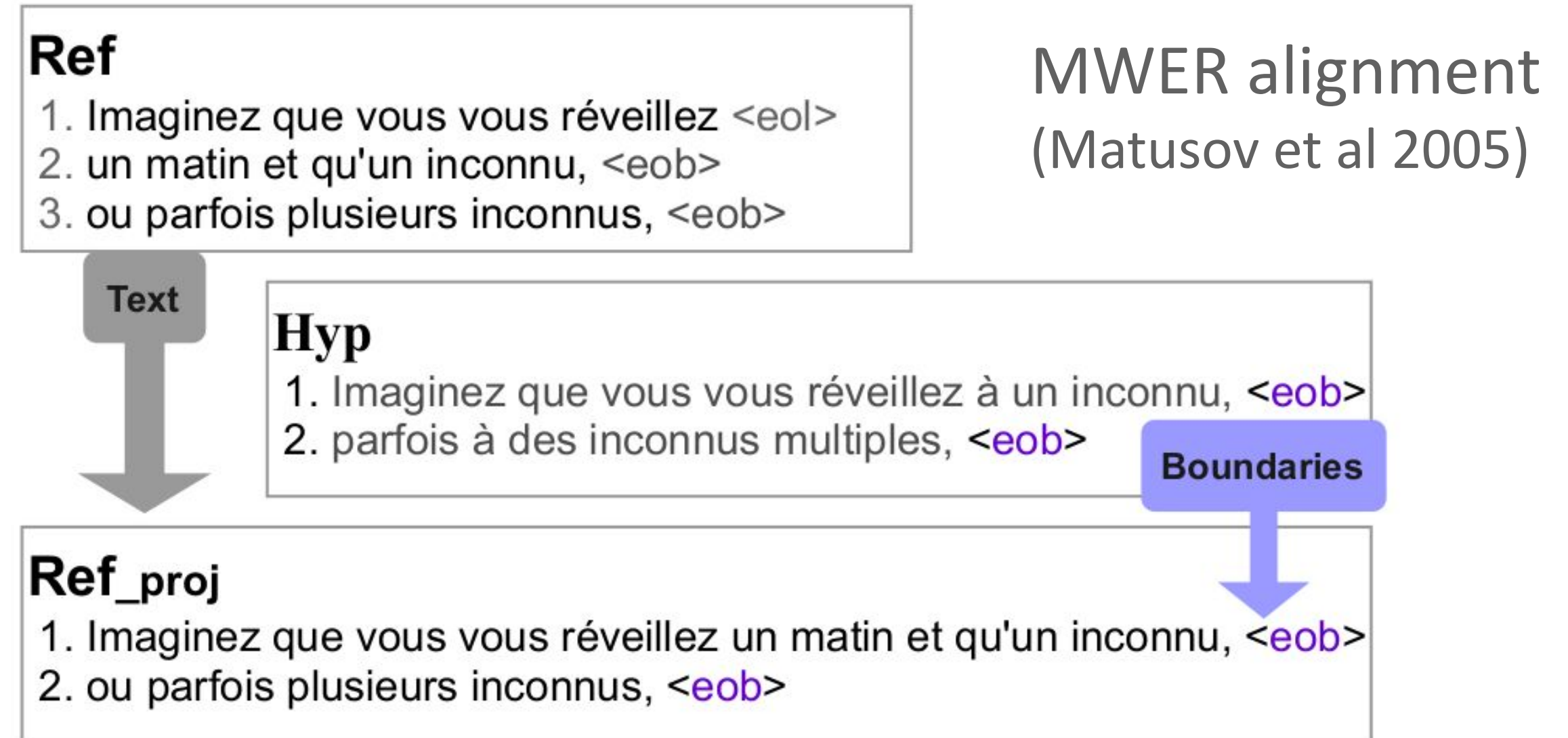
- BLEU<sub>br</sub>**
- balance among error types
  - insensitive to shift size

### Contributions

- Exp1: A **comparison of sequence segmentation metrics** for **perfect** texts
- Exp2: **Sigma**, a new segmentation score for **imperfect** texts
- Exp3: A **boundary projection method** to compute standard segmentation metrics for **imperfect** texts
- EvalSub**: A tool for computing reference-based segmentation scores for automatic subtitles

### Exp3: Boundary projection

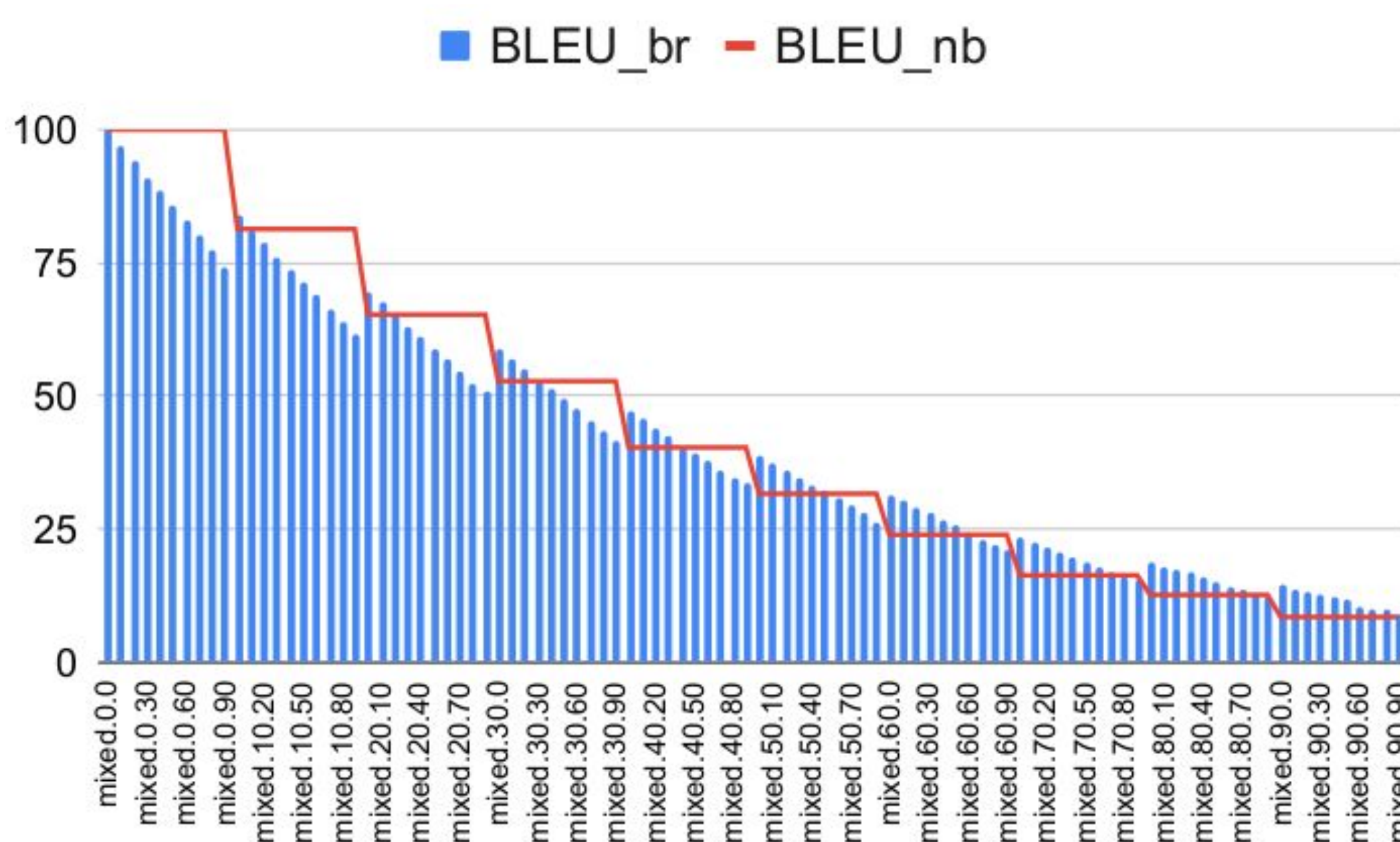
Project boundaries from Hypothesis to Reference  
Compute standard metrics



Applied on system outputs En->Fr (Karakanta et al. 2020)  
1) NMT, 2) ST cascade, 3) direct ST, 4) pretrained direct ST

### Exp2: What does BLEU<sub>br</sub> really measure?

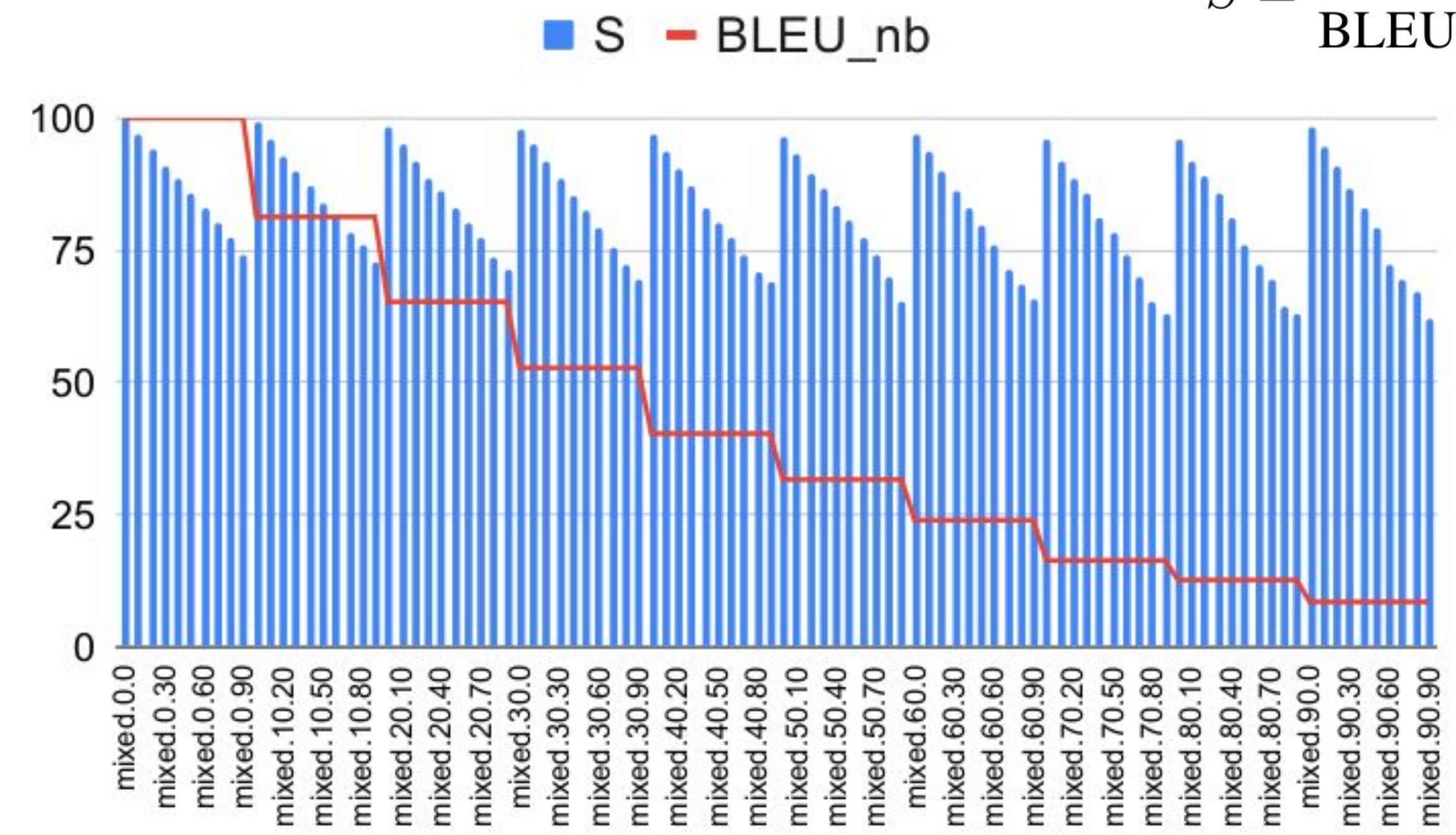
BLEU<sub>br</sub> largely depends on text quality (BLEU<sub>nb</sub>).  
The relative difference with BLEU<sub>nb</sub> cannot be used to measure segmentation quality in general.



### Sigma: a new segmentation metric

**Sigma:** ratio to a BLEU<sub>br</sub> upper bound, computed from the proportion of boundaries and n-gram precisions.  
Stable irrespective of BLEU<sub>nb</sub>

$$S = \frac{BLEU_{br}}{BLEU_{br}^+}$$



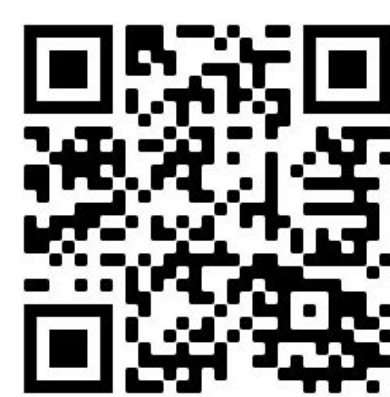
System ranking:  
NMT>Cas>e2e

System	Projected					Imperfect		S				
	P <sub>k</sub>	Windiff	SegSim	BndSim	Prec	Rec	F1		BLEU <sub>br</sub>	TER <sub>br</sub>	BLEU <sub>br</sub>	TER <sub>br</sub>
NMT	.192	.208	.979	.637	.711	.735	.723	83.18	6.87	32.16	19.38	89.2
Cas	.252	.270	.970	.519	.639	.667	.653	76.14	8.91	26.34	23.23	83.1
e2e <sub>base</sub>	.257	.277	.969	.515	.601	.667	.632	75.00	9.29	22.53	24.48	81.8
e2e <sub>pt</sub>	.252	.276	.969	.525	.610	.702	.653	74.89	9.24	26.36	23.52	81.5

Agreement on high quality output

### Data and implementation

**Data:** MuST-Cinema test set EN->FR



Code available at:  
<https://github.com/fyvo/EvalSubtitle>

### Conclusions

**Sigma** is a new BLEU-based metric suited for evaluation of subtitle segmentation in the case of end-to-end generation systems.  
Future work: correlation with human judgements