

## 8th Live Subtitling and Accessibility Symposium

# The quality of automatic and human live captions in English – and beyond

---

Pablo Romero-Fresco and Nazaret Fresno  
(GALMA, Universidade de Vigo / The University of Texas at Rio Grande Valley)



Universidade de Vigo

SWISS **TXT**



TRANS  
MEDIA  
CATALONIA

# Gion's questions

- Why is the NER model so time-consuming?
- Is it really scalable?
- Won't AI soon sweep away everything we're talking about here?



- 1 - To assess the accuracy of ASR in Galician, Basque and Catalan for live subtitling.
- 2 - To assess the accuracy of human and automatic live subtitles in English and Spanish.
- 3 - To explore the automation of the NER model
- 4 - To monitor the quality of the automatic ASR and MT tool used by the EU Parliament

UniversidadeVigo

 UNIVERSITAT  
JAUME I



The University of Texas  
Rio Grande Valley

rtve



Language	Minutes analysed	Model
Galician	475	WER and NER
Basque	225	NER
Catalan	410	WER and NER
Spanish	1000	WER and NER
English	800	WER and NER
<b>Total</b>	<b>2910</b> <b>(= 61000 subtitles)</b>	

# The quality of human and automatic subs in English

- ❑ 2018-2023
- ❑ 17.000 captions
- ❑ 800 minutes of live captions
  - { 388m human (steno and respeaking)
  - { 410m ASR
- ❑ Sky (UK), VITAC and ENCO (UK), CRTC (Canada)
- ❑ NER model

# The quality of human and automatic subs in English

Case Study	Human average	Automatic average	Year
Case Study 1 (Sky)	97% (2.5/10)	95.7% (0/10)	2018
Case Study 2 (Vitac)	98.8% (7/10)	96.3% (1/10)	2020
Case Study 3 (Enco)	99.4% (8.5/10)	96.3% (1.5/10)	2021
Case Study 4 (Canada)	98.8% (7/10)	97.8% (4.5/10)	2020-2022
<b>Totals</b>	<b>98.5% (6.5/10)</b>	<b>96.5% (1/10)</b>	

# The quality of human and automatic subs in English

Jim Pattison, 2021-2022: 98.3% (31 samples)

Miracle Channel, 2021-2022: 98.3% (24 samples)

Case Study 4 (Canada)	98.8% (7/10)	97.8% (4.5/10)	2020-2022
-----------------------	--------------	----------------	-----------

# The quality of human and automatic subs in English

## Automatic captions:

- Close to 98%
- Slightly over 98% with human assistance
- Paves the way for other languages
- Still issues with punctuation, proper nouns, numbers, speaker ID
- Always (near) verbatim, so often fast and difficult to read



# The quality of human and automatic subs in English

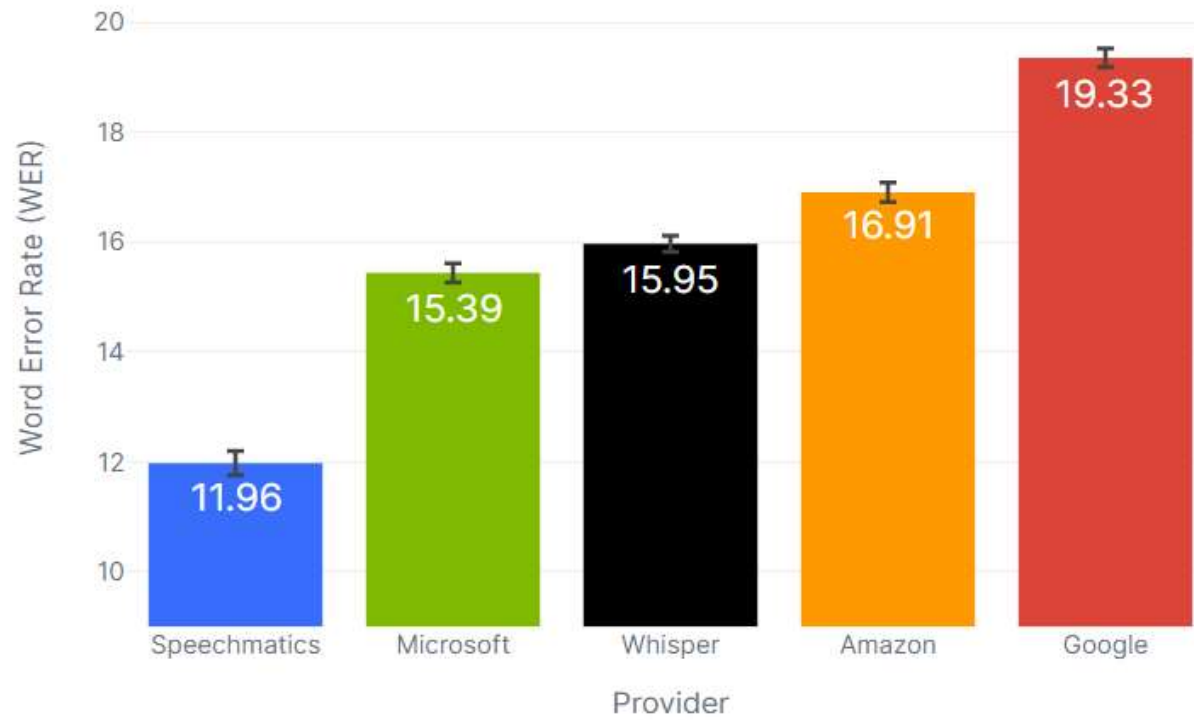
Human captions:

- Excellent (steno) or very good (respeaking)
- More accessible: fewer errors and more readable speeds

# The quality of human and automatic subs in English



# The WER War



# The WER War

\*\* We are aware of the limitations of WER, one major issue being that errors involving misinformation are given the same weight as simple spelling mistakes. To address this, we normalize our transcriptions to reduce penalties for differences in contractions or spelling between British and American English that humans would still consider correct. Going forward, we intend to adopt a metric based on NER.

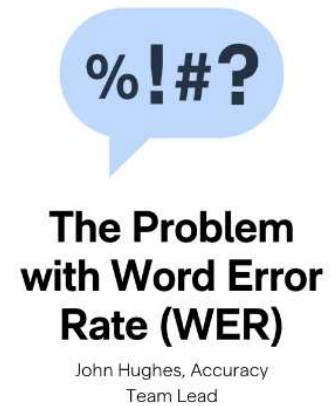
# The WER War

Blog - Technical

Oct 14, 2022 | Read time 4 min

## The Problem with Word Error Rate (WER)

In part one of a two-part piece, Accuracy Team Lead, John Hughes, explains how and why Word Error Rate as an accuracy measure is outdated and often misaligned with human judgment. Part Two looks at where we go next.



John Hughes  
Accuracy Team Lead



# WER vs NER

	Speechmatics		Whisper		Pre-AI ASR	
English speech						
English interview						
Spanish speech						
Spanish interview						

# WER vs NER

	Speechmatics		Whisper		Pre-AI ASR	
English speech	3.6%		3.8%		3%	
English interview	15.7%		19.1%		25.7%	
Spanish speech	1.7%		2%		5.2%	
Spanish interview	14%		10.3%		27.6%	







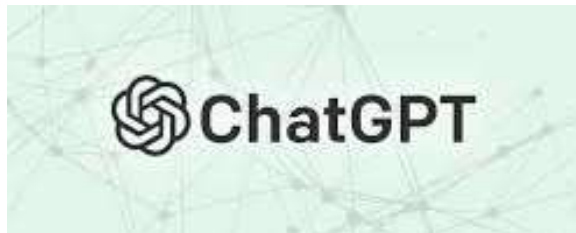
# WER vs NER

	Speechmatics		Whisper		Pre-AI ASR	
English speech		99.6%		99.7%		98.9%
English interview		99.4%		99.9%		97%
Spanish speech		99.2%		99.9%		97.1%
Spanish interview		99.1%		99.7%		95%

# WER vs NER

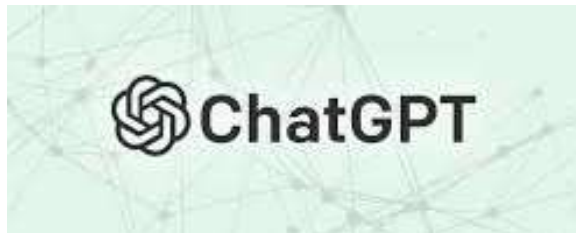


# Can NER be automated?



- 300 examples of NER
- Distinction between minor, standard and serious errors
- 3 months
- 2 iterations

# Can NER be automated?



Confusion matrix:

GPT Score	Human Label					All
	0.25	0.5	0.75	1	1.25	
0.25	79	26	0	6	1	112
0.5	39	24	3	4	0	70
0.75	11	3	1	0	0	15
1	33	41	4	14	0	92
1.25	4	3	2	1	0	10
1.5	5	1	0	0	0	6
1.75	1	2	0	0	0	3
2	3	2	1	1	0	7
2.25	1	0	0	0	0	1
All	176	102	11	26	1	316

# Can NER be automated?



- 40% success rate
- Good with minor errors
- Sometimes struggles to differentiate standard errors (meaning is lost) from serious errors (meaning is changed)

## Conclusions (2017-2022)

- ❑ 2017- late 2022: improvement of autosubs from unusable to useable
- ❑ Human captions still better and more accessible (edited)

# Conclusions (2023)

New AI-powered ASR: potentially more accurate than human\*

\* verbatim, no or worse speaker ID

Important decisions adopted on the basis of unreliable WERs

Potential automation of NER

Future of intra and interlingual human and automatic subs?



# Gion's questions

- Why is the NER model so time-consuming?
- Is it really scalable?
- Won't AI soon sweep away everything we're talking about here?

# 8th Live Subtitling and Accessibility Symposium

## The quality of automatic and human live captions in English – and beyond

---

Pablo Romero-Fresco and Nazaret Fresno  
(GALMA, Universidade de Vigo /  
The University of Texas at Rio Grande Valley)



Universidade de Vigo

**SWISS** **TXT**



TRANS  
MEDIA  
CATALONIA