

完璧でない機械学習システムによる／のための 行動変容インタラクション

Hiromu Yakura

University of Tsukuba, Japan

2023年7月25日 @ 第67回AIセミナー

Self-Introduction

- **Hiromu Yakura**
 - Ph.D. student at University of Tsukuba
 - Satellite lab at Media Interaction Group, AIST Tsukuba
 - Google / Microsoft Research Ph.D. Fellow

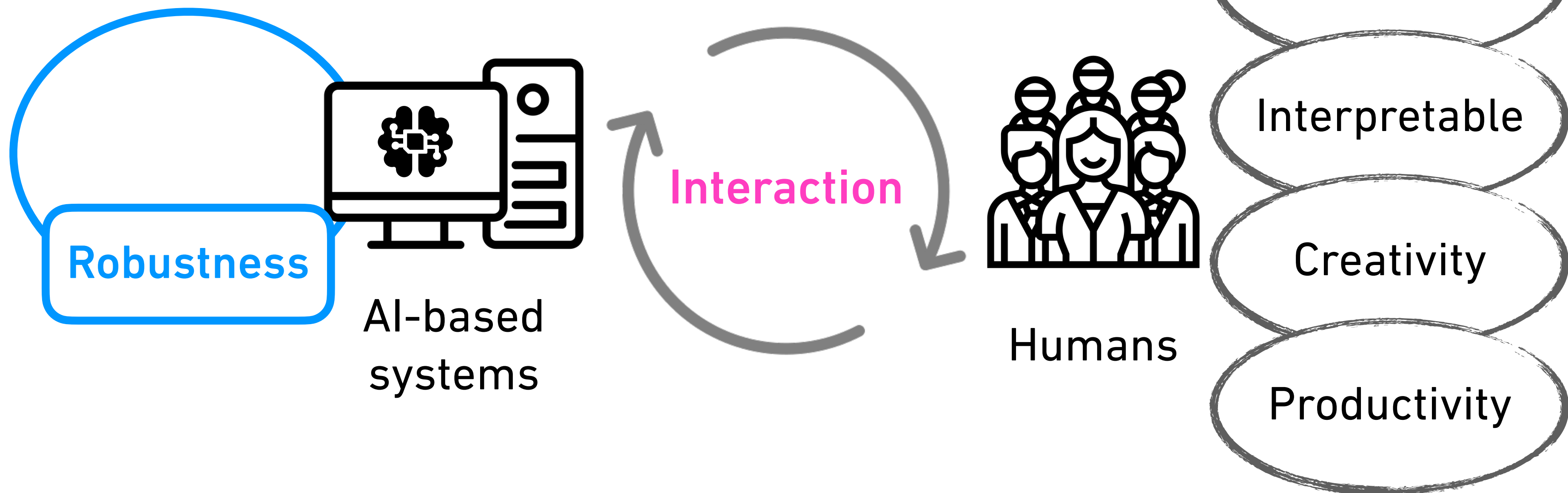


@hiromu1996

Research context

- Specific research interest: **ML** + **HCI**
 - How to apply **machine learning** in a **human-centric** manner?

ML techniques



Robust Audio Adversarial Example for a Physical Attack

Hiromu Yakura^{*†}, Jun Sakuma^{*†}

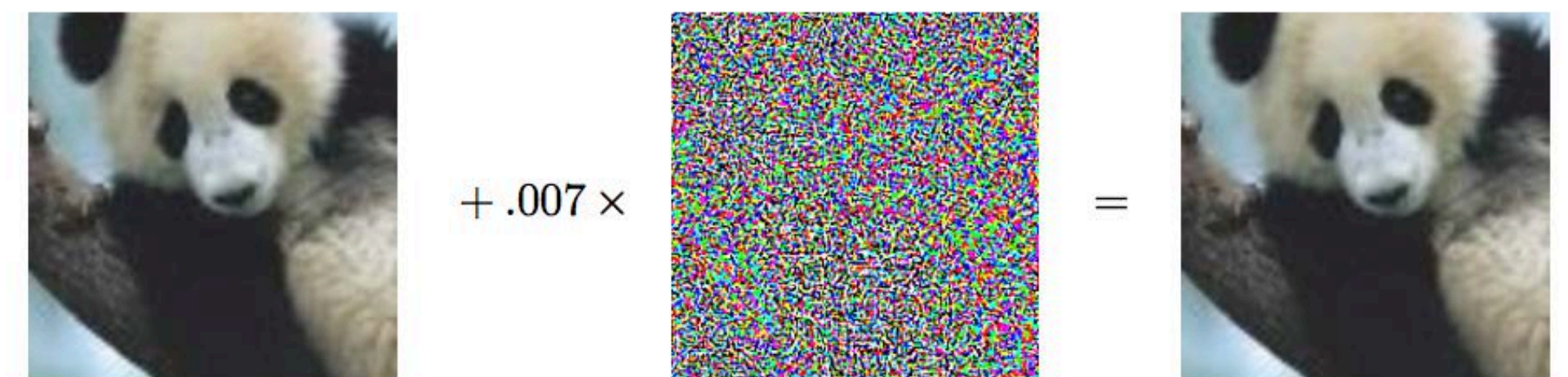
^{*} University of Tsukuba, Japan

[†] RIKEN Center for Artificial Intelligence Project, Japan

IJCAI 2019

Background: Risk of abusing AI using adversarial example

- It is known that **we can mislead ML models** by intentionally adding a small noise to the inputs.
- In this case, the right image is classified as gibbon.



x
“panda”
57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

$=$

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

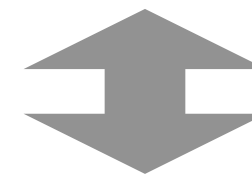
Using this adversarial example, **we can abuse AI-based systems** without being noticed by humans.

Background: Risk of abusing AI using adversarial example

- The risk of such an attack is not limited to images, but also to speech recognition.
- In particular, speech recognition is widely used in the form such as Siri or Google Home.



- If we can mislead speech recognition models to transcribe specific words, such **voice assistants can also be abused.**



However, abusing speech recognition in the real world has been considered difficult due to noise. [Carlini+, 2018]

Proposal: Robust attack for speech recognition

**We realized such an attack in the real world
by simulating reverberations and noises in Tensorflow.**

Original audio



Manipulated audio transcribed as

"hello world"



"open the door"



- In a listening experiment involving 50 participants, we confirmed no one could tell these hidden messages.

Generate (non-software) Bugs to Fool Classifiers

Hiromu Yakura^{*†}, Youhei Akimoto^{*†}, Jun Sakuma^{*†}

^{*} University of Tsukuba, Japan

[†] RIKEN Center for Artificial Intelligence Project, Japan

AAAI 2020

Proposal: Attacking self-driving cars with moth-like stickers

- This mechanism is also applicable to deceive self-driving cars.
- What if the cars recognize a STOP sign as Speed 80?
- This example can cause such a mistake but looks too suspicious not to be noticed by humans.



We showed that **these moth-like stickers can mislead ML models without making humans feel suspicious.**

Interaction Design to Leverage Fallible Machine Learning Models

Basic idea: How to overcome the fallible of AI models

No matter how much technical improvements we make,
AI-based systems will make mistakes.

- AI models just perform inference based on the trends they found in the data given in advance.
- Thus, we can't deny the possibility of their mistakes, even if we don't intentionally attack them.



**Design an interaction in which humans and AI-based systems
can collaborate effectively even when AIs are not perfect**

Mindless Attractor: A False-Positive Resistant Intervention for Drawing Attention Using Auditory Perturbation

Riku Arakawa[†] and Hiromu Yakura[‡]
(equal contribution)

[†] The University of Tokyo, Japan

[‡] University of Tsukuba, Japan

ACM CHI 2021

Background: Limitation of alerting intervention



[Gupta+, '16]

- ML models can detect the moment when people are not engaging.
- It is possible to alert distracted students in video lectures.
- But, **misinformed alerts caused by false positives disrupt the students**, which leads them to distrust the system.
- In addition, **alerting like "keep focused" would not help them keep focused** unless they have strong motivation.



- What is the best way to intervene in distracted students while we cannot deny the possibility of false positives?

Demo: Mindless Attractor



Proposal: Mindless Attractor

Humans often intentionally or unconsciously change the volume or pitch of speech to draw listeners' attention.

- Our brains are known to respond to such signals. [Zatorre+, '07]



- It computationally changes the volume or pitch for a moment to draw attention without consuming conscious awareness.
- This is machine learning-friendly because it won't be frustrating even when activated by false-positive detection.
- In our study, we confirmed its effectiveness to help refocus even when the students are unconscious of the changes.

BeParrot: Efficient Interface for Transcribing Unclear Speech via Respeaking



Riku Arakawa[†]

Carnegie Mellon University



Hiromu Yakura[†]

University of Tsukuba

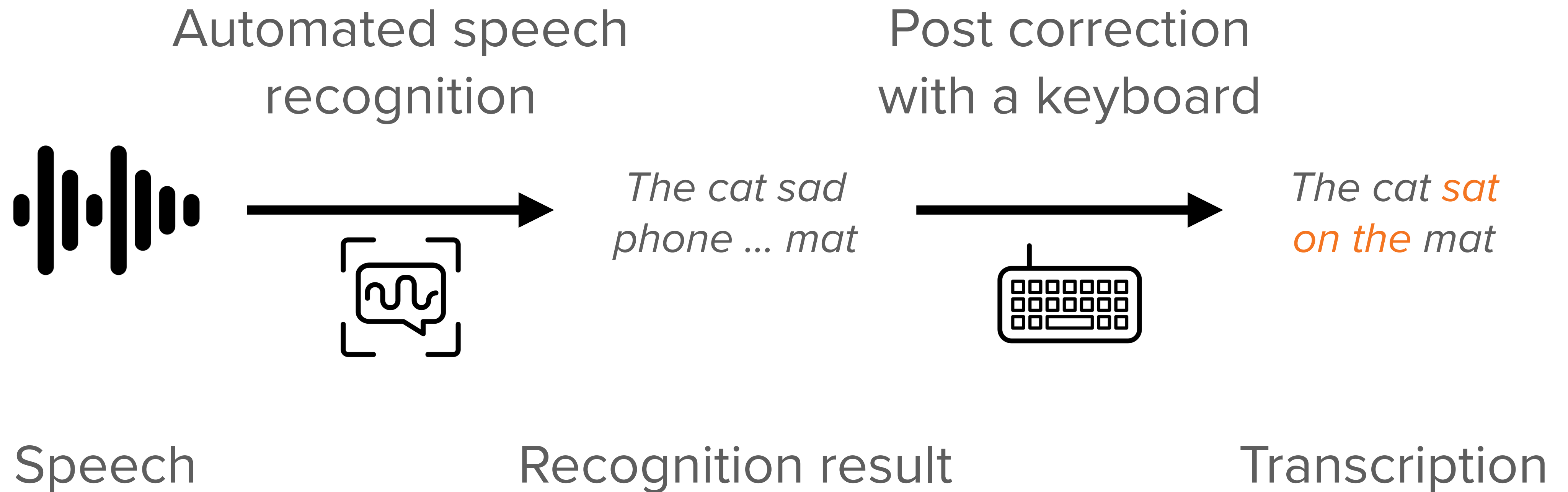


Masataka Goto

National Institute of Advanced Industrial
Science and Technology (AIST)

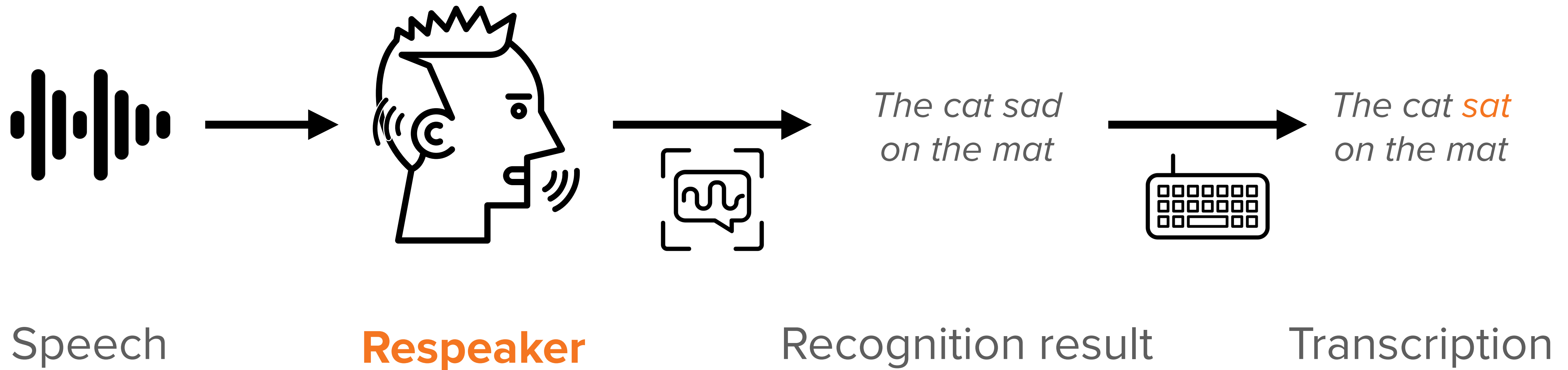
[†] Equal contribution

How to transcribe speech with recognition models



What if the speech is unclear?

Respeaking



Respeaking can reduce the number of error corrections

0:23 / 3:00

再生速度 x1.0
読み上げ長さ 10.0秒

次へ
やり直す

自動認識結果

前： 現在地
今： 現在永遠の地下鉄の文学のことではの古典の
次： 世界みたいなのが入った

読み上げ結果

確定済み

ただいまご紹介に預かりました大変良いところで素晴らしい環境だと思います今日は文学の現在 というテーマでお話するわけですが

This is a demo where a Japanese user tries to transcribe a historical speech.

Respeaking in Practice: Live Captioning

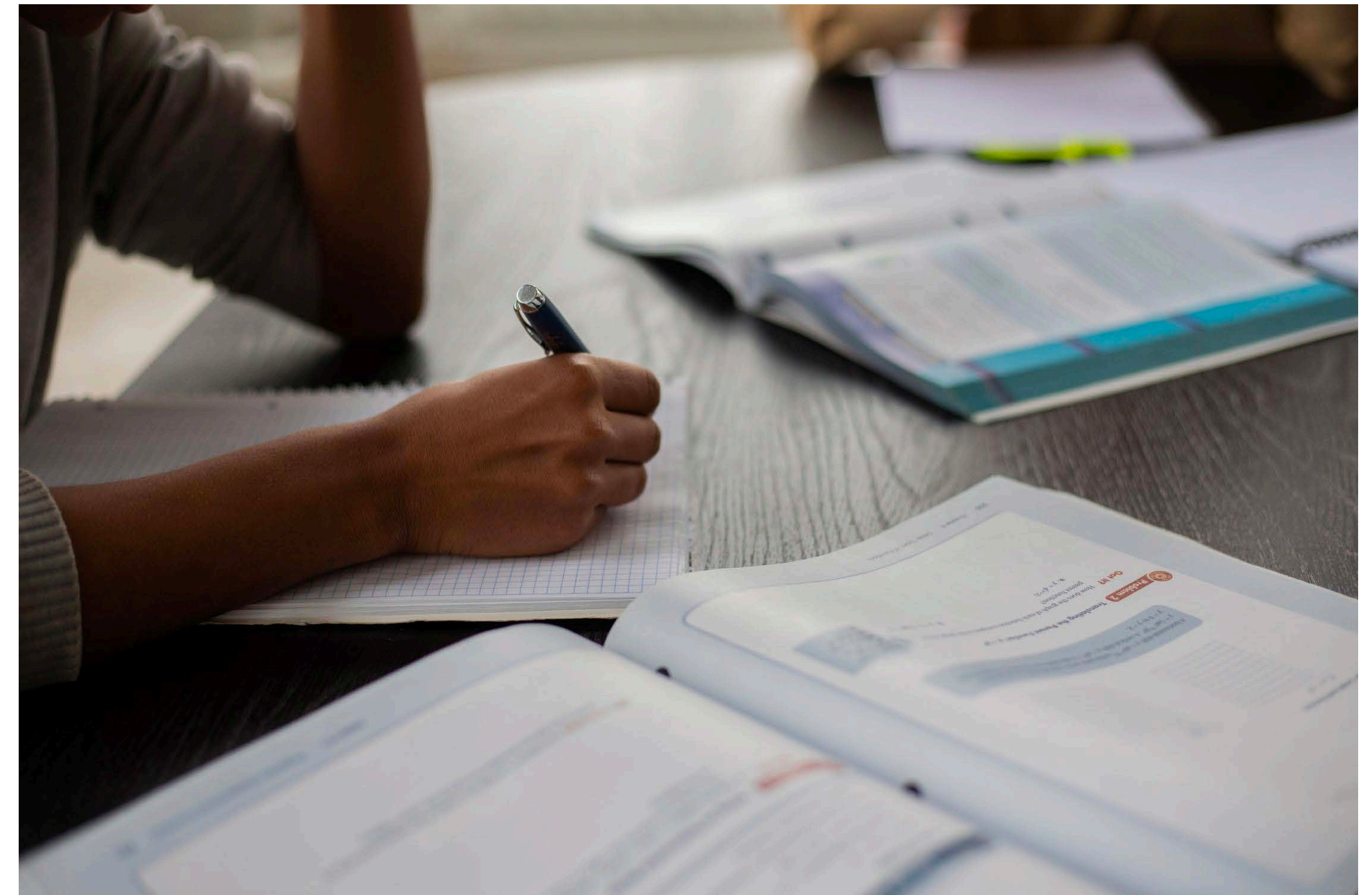


Source: <https://www.nhk.or.jp/str/publica/rd/182/3.html>

Hurdle of respeaking

- Respeakers should be able to repeat the provided speech clearly without stuttering or stammering
- Respeakers are required not only to repeat the speech but also to memorize the speech content

To be a professional respeaker ...



There is a training program of 75 hours

Proposed System: BeParrot

The screenshot displays the BeParrot interface with the following components:

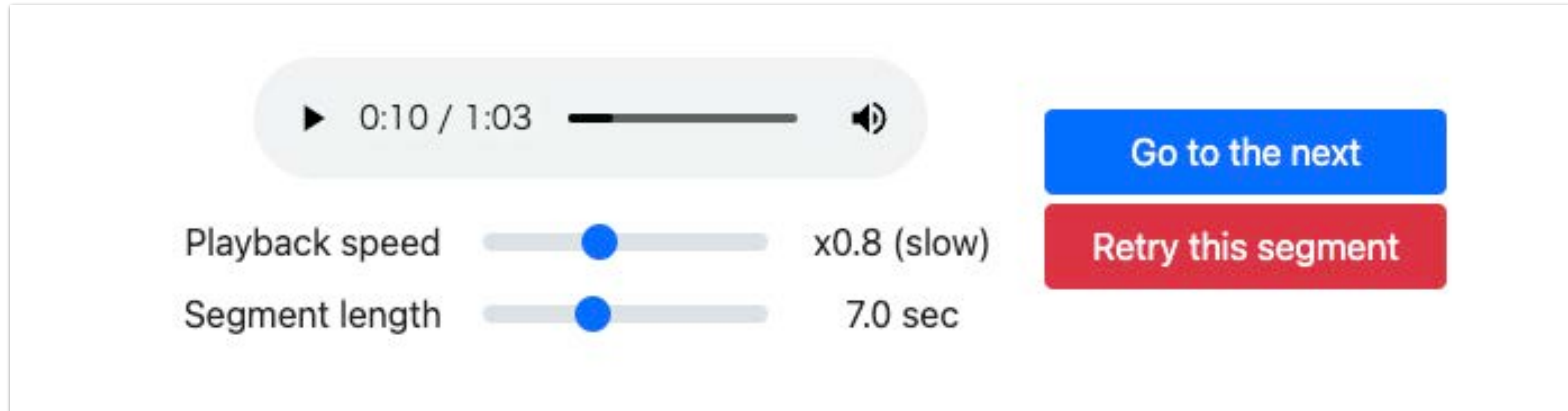
- Audio Controls:** A progress bar at 0:50 / 1:03, a volume icon, a "Go to the next" button, and a "Retry this segment" button. Below these are sliders for "Playback speed" (set to x0.9) and "Segment length" (set to 7.0 sec).
- Pre-recognized content (for reference):** A dark gray box containing three lines of text:
 - Prev: we will be able to speed up at 10
 - Curr: but I will
 - Next: be able to bring in the words of the old Negro spiritual free
- Recognition result:** A bright blue box containing the text "but I will".
- Final transcription:** A dark gray box containing the text:

I have a dream my poor little children one one day live in the nation
where they will not be judged by the color of their skin but by the
content of a character I have a dream today will be able to speed up at
the
- Words misrecognized:** A section at the bottom left listing:
 - live: labor (1 time)
 - the: 10 (1 time)
- Help Link:** A yellow button at the bottom right labeled "Having trouble?" with a question mark icon.

Key features:

- **Parameter adjustment**
- **Pronunciation feedback**

System Feature: Parameter Adjustment



Slower playback speed → Help avoiding stuttering or stammering

Shorter segment length → Reduce the demand of memorizing content

System Feature: Pronunciation Feedback

Words not recognized

at: 1 time

Words misrecognized

the: 10, day (2 times)

live: labor (1 time)

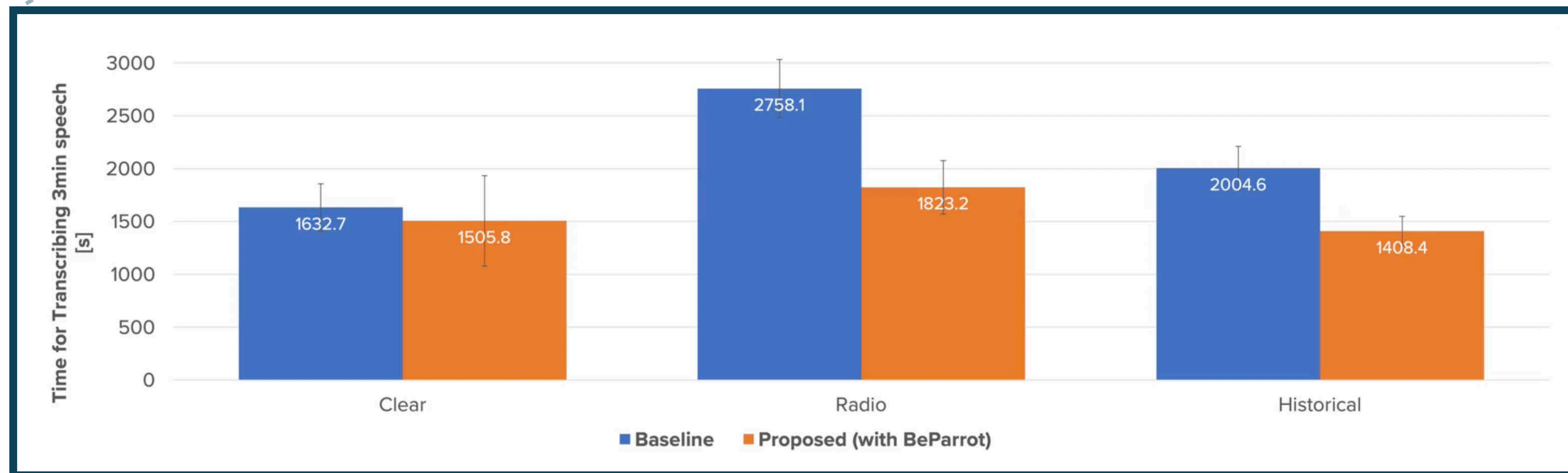
bring: ring (1 time)

Automatically calculated from the correction history

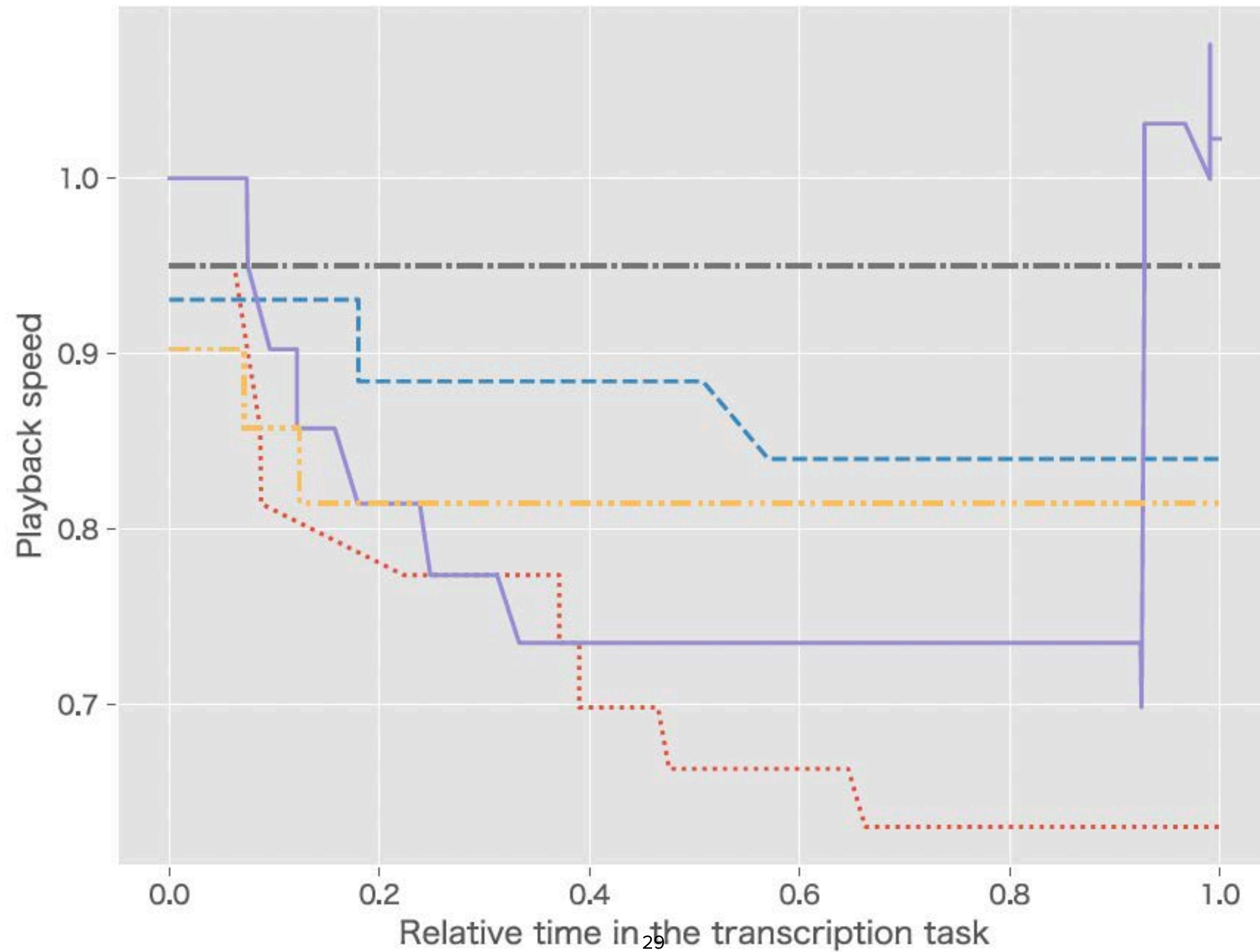
Study Results: Time

Speech type	Time			CER (%)		
	Baseline (s)	Proposed (s)	Reduction (%)	ASR	Baseline	Proposed
Clear	1632.7 (± 224.6)	1505.8 (± 426.8)	7.8	6.16	3.73 (± 0.57)	5.75 (± 1.02)
Radio	2758.1 (± 273.6)	1823.2 (± 253.4)	33.9	30.72	19.05 (± 1.56)	24.81 (± 2.69)
Historical	2004.6 (± 205.0)	1408.4 (± 140.4)	29.7	48.18	29.25 (± 5.16)	29.10 (± 2.17)

32.1 % time reduction for transcribing
unclear speech (radio + historical speech)



Study Results: Speed History



Conclusion

- HCI can provide a glue between users and fallible ML systems.
- Better interaction design can **reduce the risk of false positives** in inducing behavior change of users.
- It is also possible to **train users to behave ML friendly** to foster effective human-AI collaboration.