



Human Interactivity
and Language Lab



[Speakers and tutors](#) [Programme](#) [Practical info](#) [Application](#) [Organizers](#)

September 18-23, 2023, Poland

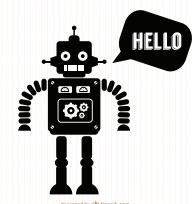
4th Summer School on Social Human-Robot Interaction



1

© 2023 The University of Sheffield

Should a Robot Speak? *(if so, why, when and how)*



Prof. Roger K. Moore

Chair of Spoken Language Processing
Dept. Computer Science, University of Sheffield
(Visiting Prof., Language Sciences, University College London)
(Visiting Prof., Bristol Robotics Lab.)



Sheffield Robotics Seminar

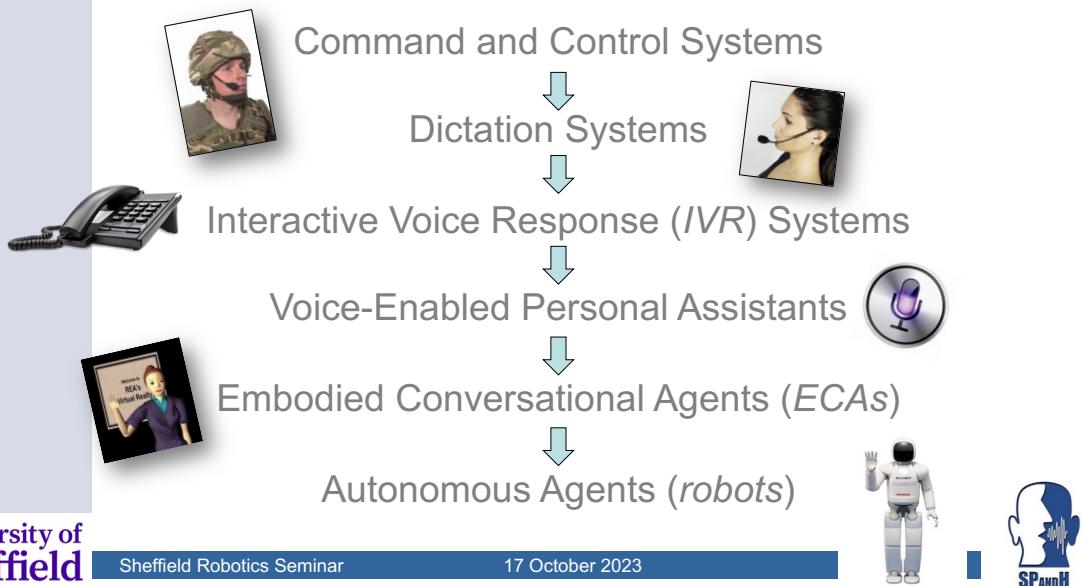
17 October 2023

slide 2

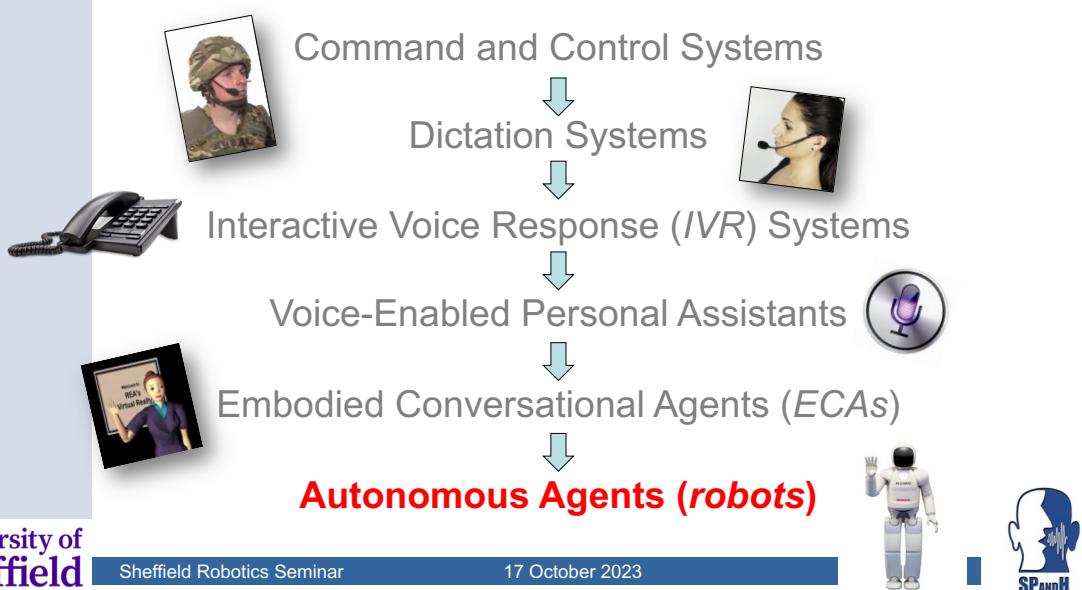
2

1

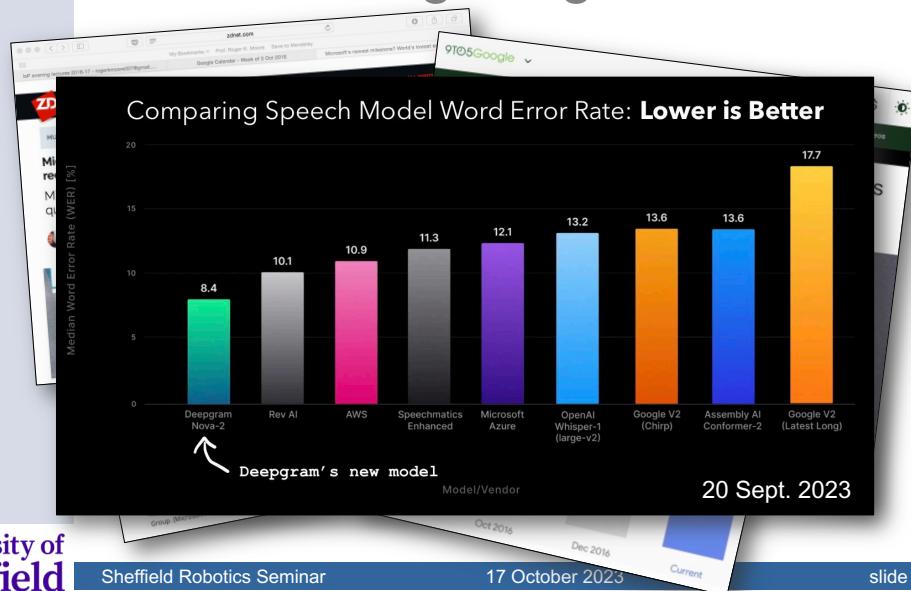
Talking with Machines



Talking with Machines



Amazing Progress



Sheffield Robotics Seminar

17 October 2023

slide 5



5

Where Are We Going?

GLOBAL VOICE ASSISTANT MARKET



North America
Largest Market
By Region (2019E)



APAC
Fastest-Growing Market
By Region (2020-2030)

2019
Market Size
\$1,723.6
million

2030
Market Size
\$26,872.6
million

Market
Growth Rate
(2020-2030)
29.7 %



Sheffield Robotics Seminar

17 October 2023

slide 6

6

© 2023 The University of Sheffield

Where Are We Going?

GLOBAL

North Largest By Reg

BUSINESS INSIDER UK

TECH

The next generation of Siri-like assistants will be robots living in your home

■ ANTONIO VILLAS-BOAS JUN. 22, 2015, 8:40 PM 419

FACEBOOK IN LINKEDIN TWITTER

Most robots so far have remained on testing platforms or on stages at showcase events to show off a company's technological ability.

The Pepper robot built by Japanese companies Alderbaran and Softbank, on the other hand, can be bought online and used at home right now. Except you may have to wait a long time as the first batch of 1000 Peppers was sold out in about a minute, and it's only available in Japan.

Pepper, the personal digital assistant robot.

Business Insider

Pepper, the personal digital assistant robot.

Market Growth Rate (2020-2030) 29.7 %

https://www.psmarketing.com/-analysis/voice-assistant-market

Sheffield Robotics Seminar 17 October 2023 slide 7 SPANDH

7

© 2023 The University of Sheffield

Recent Trends

Introducing ChatGPT

We've trained a model called ChatGPT which interacts in a conversational way. This unique format makes it possible for ChatGPT to answer questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

By OpenAI Read more about ChatGPT

OpenAI

28 Aug 2022 1 Jan 2023 25 Sept 2023

Coffee

Alexa

ChatGPT

Research API ChatGPT Safety Company

ChatGPT can now see, hear, and speak

We are beginning to roll out new voice and image capabilities in ChatGPT. They offer a new, more intuitive type of interface by allowing you to have a voice conversation or show ChatGPT what you're talking about.

University of Sheffield

Sheffield Robotics Seminar 17 October 2023 slide 8 SPANDH

8

© 2023 The University of Sheffield

The ‘State-of-the-Art’

- There is steady year-on-year progress
- Improvements come from:
 - increase in available computer power
 - latest machine-learning paradigms
 - huge real-world training corpora
 - open tools / benchmark testing / challenges
- Progress has not come about as a result of deep insights into human spoken language
- Spoken language technology is ...
 - fragile (*in ‘real’ conditions*)
 - expensive (*to port to new applications / languages*)
 - inefficient (*trained on more data than a person hears in a lifetime*)
 - biased (*trained on un-curated data*)
 - untrustworthy (*e.g. hallucinates misinformation*)
 - ecologically damaging (*model training has a significant carbon footprint*)
- It is easy to underestimate the richness and complexity of spoken language interaction
- It is not ‘natural’ to talk to a machine!
- The availability of open tools and data is de-skilling



University of Sheffield

Sheffield Robotics Seminar 17 October 2023 slide 9

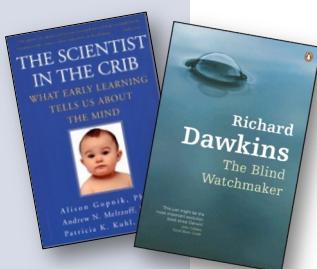


9

© 2023 The University of Sheffield

Speech is not just another Modality

“Spoken language is the most sophisticated behaviour of the most complex organism in the known universe!”



Simply interfacing state-of-the-art speech technology with a state-of-the-art robot does not lead to effective human-robot interaction (*talking and listening is not enough*)



University of Sheffield

Sheffield Robotics Seminar 17 October 2023 slide 10



10

Speech is not just Audible Text



Sheffield Robotics Seminar

17 October 2023

slide 11



11

Speech is ...

- variable
- ambiguous
- effortful
- contrastive
- prosodic
- adaptive
- context-dependent
- meaningful
- referential
- indexical
- rhetorical
- paralinguistic
- personalised
- affective
- multimodal
- contaminated

Prediction and Entropy of Printed English
By C. E. SHANNON
(Manuscript Received Sept. 15, 1950)

A new method of estimating the entropy and redundancy of a language is described. This method exploits the knowledge that language statistics provide by themselves who speak the language. The redundancy, which depends on experiments, is given for the next letter after the preceding text of 10 letters. Results of experiments in prediction are given, and some properties of an ideal predictor are developed.

1. INTRODUCTION

In a previous paper¹ the entropy and redundancy of a language have been defined. The entropy is a statistical parameter which measures, in a certain sense, how much information is produced on the average for each letter of a text in the language. If the language is translated into binary digits (0 or 1) in the most efficient way, the entropy H is the average number of binary digits required per letter of the original language. The redundancy, on the other hand, measures the amount of constraint imposed on a text by the language due to its statistical structure, e.g., in English the high frequency of the letter E , the strong tendency of H to follow T or of T to follow E , etc. It was estimated that when statistical effects extend over 10 letters, the entropy is roughly 2.3 bits per letter, and that eight letters are considered the redundancy is about 30 per cent.

Since then a new method has been found for estimating these quantities, which is more sensitive and takes account of long range statistics, influenced by extending over phonetic sentences, etc. This method is based on a study of the predictability of English; how well can the next letter of a text be predicted when the preceding N letters are known. The results of some experiments will be given, and a theoretical analysis of some of the properties of ideal prediction. By combining the experimental and theoretical results it is possible to estimate upper and lower bounds for the entropy and redundancy. From this analysis it appears that, for ordinary literary statistical effects (up to 100 letters) reduce the redundancy by about 8 bps when compared with a letter, with a corresponding increase in the entropy of about 39 bps.

Shannon, C. E. (1951). Prediction and entropy of printed English. *The Bell System Technical Journal*, 50–64.

The figure contains two side-by-side line graphs, one for 'SR' (Self-Reducing) and one for 'IR' (Information Rate). Both graphs plot 'Value (SR or IR)' on the y-axis against 'Language' on the x-axis. The x-axis lists languages: All, JPN, SPA, EUS, FIN, ITA, SRP, SOR, CAT, TUR, FRA, ENG, DEU, HUN, CMN, YUE, VIE, THA. The y-axis ranges from 4 to 60 for SR and 30 to 60 for IR. Each language has three overlapping bell-shaped curves representing different language families: Austroasiatic (pink), Indo-European (orange), and Sino-Tibetan (yellow). A red arrow points to the 'IR' graph for 'ITA' with the label '~8 bps'. Another red arrow points to the 'IR' graph for 'All' with the label '39 bps'.

Coupé, C., Oh, Y., Dedić, D., & Pellegrino, F. (2019). Different languages, similar encoding efficiency: Comparable information rates across the human communicative niche. *Science Advances*, 5(9).

Sheffield Robotics Seminar

17 October 2023

slide 12

12

6

The Information Rate of Spoken Language

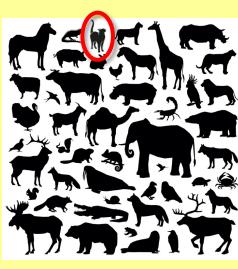
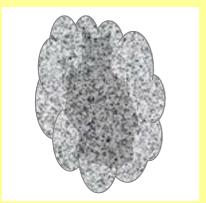


Moore, R. K. (2023). Pragmatics, synchronics and energetics in spoken language – an information theoretic perspective. In *Workshop on Limits and Benefits of Information-Theoretic Perspectives in Spoken Communication* (LBIT), Dublin.

- Of course, spoken language is *not* a fixed code with a constant information rate
- What people say is conditioned on critical *causal* factors such as ...
 - their situated and embodied circumstances: '**pragmatics**'
 - the temporal evolution of events: '**synchronics**'
 - the level of effort that they are prepared to devote to their behaviour: '**energetics**'
- In other words, the information rate in spoken language varies as a function of the pragmatic, synchronic and energetic context



Pragmatics: situational context

<p>Entropy = 5.58 bits</p>  <p>"Ooh, look at that cat"</p> <p>Communicated = 5.58 bits</p>	<p>Entropy = 1 bit</p>  <p>"Ooh, look at that cat ... the one on the left"</p> <p>Communicated = 1 bit</p>	<p>Entropy = 1 bit</p>  <p>"Ooh, look at that cat ... the one in front and slightly to the left"</p> <p>Communicated = ~80 bits</p>	<p>Entropy = 1 bit</p>  <p>"Ooh, look at that ... I think it's a cat, or maybe two cats with one sitting in front of the other as I can see what looks like three ears with one head lower than the other and slightly to the left?"</p> <p>Communicated = ~440 bits</p>
---	---	---	---



Energetics: *motivational context*

Transmitted Information \propto Effort



"Why do I have to shout?"

"Where's the newspaper?"

"I don't know"

Hawkins, S. (2003). Roles and representations of systematic fine phonetic detail in speech understanding. *Journal of Phonetics*, 31, 373-405.

"Huh? Speak more clearly!"



Sheffield Robotics Seminar

Silhouette images from Freepik

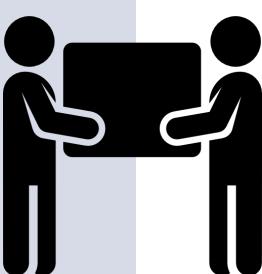
17 October 2023

slide 15



15

Synchronics: *temporal context*



"To me" ... "To you"



"Look at that ..."



"Three, two, one ... GO!"



Sheffield Robotics Seminar

Silhouette images designed by kjpargeter / Freepik

17 October 2023

slide 16



16

What is Language Like?



Cummins, F. (2011). Periodic and aperiodic synchronization in skilled action. *Frontiers in Human Neuroscience*, 5(170), 1–9.



University of
Sheffield

Sheffield Robotics Seminar

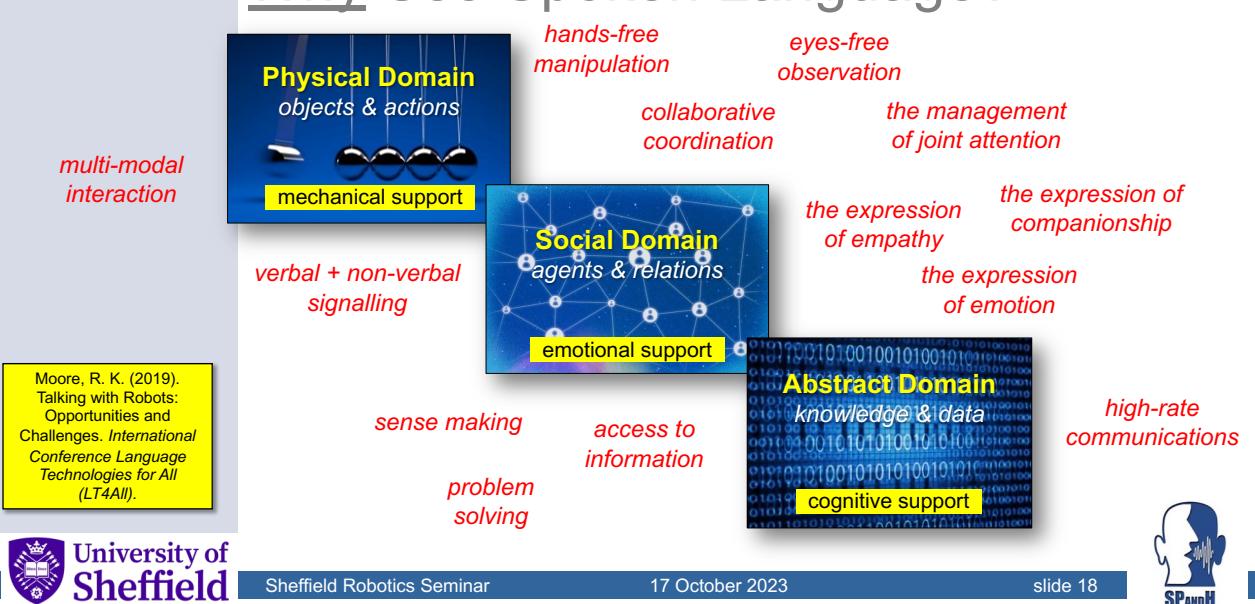
17 October 2023

slide 17



17

Why Use Spoken Language?



University of
Sheffield

Sheffield Robotics Seminar

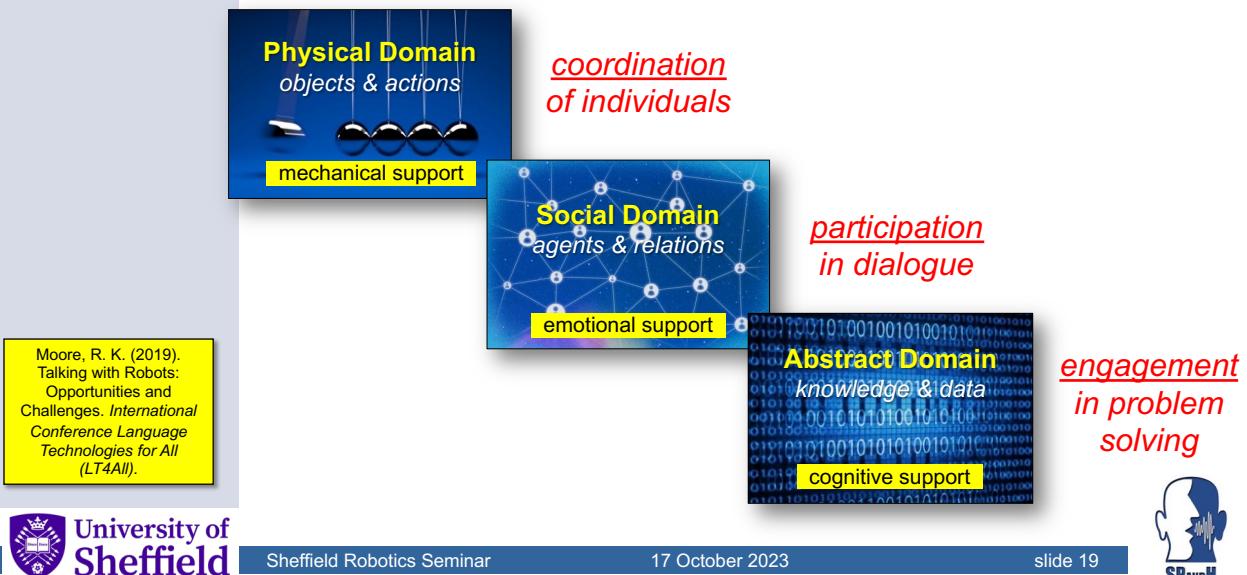
17 October 2023

slide 18

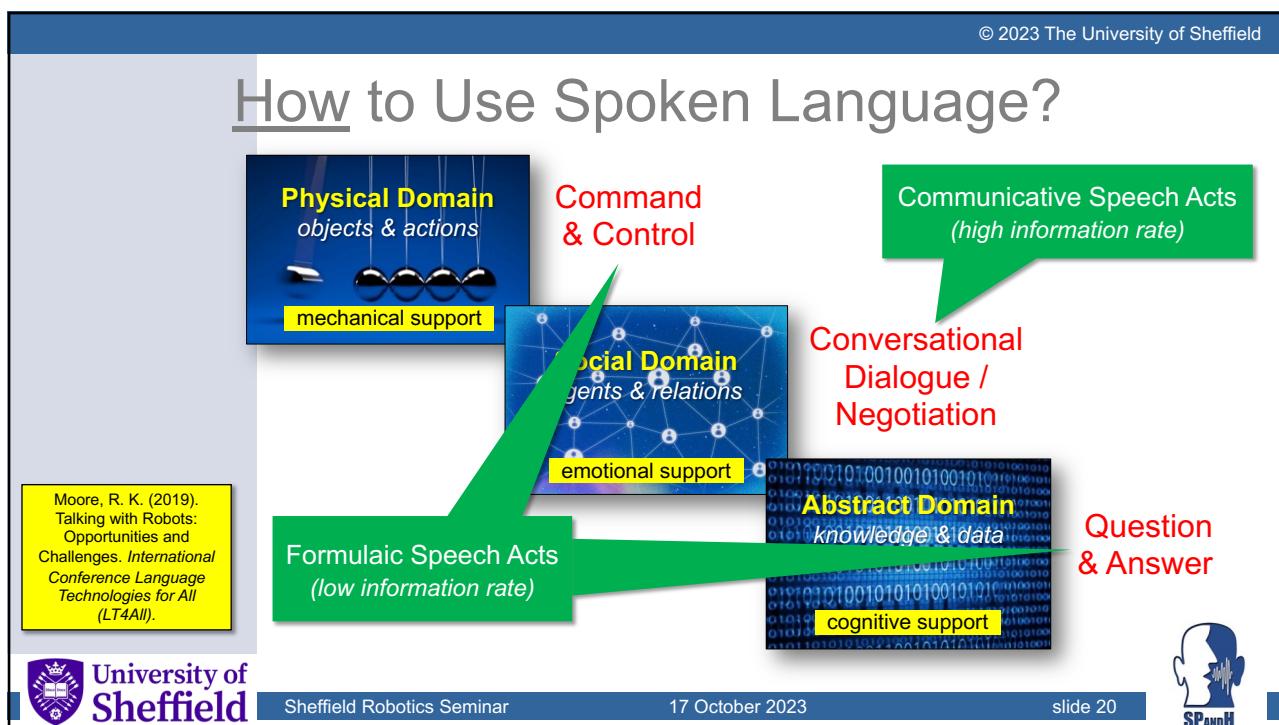


18

When to Use Spoken Language?



19



20

© 2023 The University of Sheffield

The ‘State-of-the-Art’

- There is steady year-on-year progress
- Improvements come from:
 - increase in available computer power
 - latest machine-learning paradigms
 - huge real-world training corpora
 - open tools / benchmark testing / challenges
- Progress has not come about as a result of deep insights into human spoken language
- Spoken language technology is ...
 - fragile (*in ‘real’ conditions*)
 - expensive (*to port to new applications / languages*)
 - inefficient (*trained on more data than a person hears in a lifetime*)
 - biased (*trained on un-curated data*)
 - untrustworthy (*e.g. hallucinates misinformation*)
 - ecologically damaging (*model training has a significant carbon footprint*)
- It is easy to underestimate the richness and complexity of spoken language interaction
- It is not ‘natural’ to talk to a machine!
- The availability of open tools and data is de-skilling



University of Sheffield

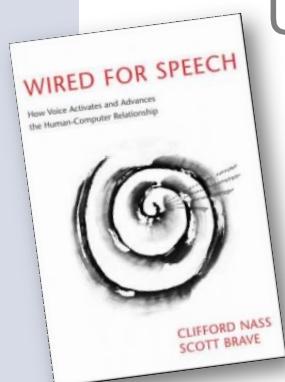
Sheffield Robotics Seminar 17 October 2023 slide 21



21

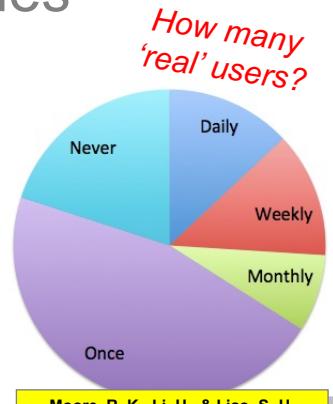
© 2023 The University of Sheffield

Usability Issues



“Voice interfaces have become notorious for fostering frustration and failure.”

- noise
- accents
- understanding
- privacy
- embarrassment
- alternative GUIs
- task familiarity
- limited functionality



How many ‘real’ users?

Frequency	Percentage
Never	~25%
Daily	~20%
Weekly	~25%
Monthly	~10%
Once	~15%

Moore, R. K., Li, H., & Liao, S.-H. (2016). Progress and prospects for spoken language technology: what ordinary people think. *INTERSPEECH* (pp. 3007–3011). San Francisco, CA.

University of Sheffield

Sheffield Robotics Seminar 17 October 2023 slide 22



22

Alexa is not a Robot!

(and a robot may not be a personal assistant)



"Alexa, play morning playlist."



Hey Siri



"A robot is an actuated mechanism programmable in two or more axes with a degree of autonomy, moving within its environment, to perform intended tasks."

<http://www.leorobotics.nl/definition-robots-and-robotics>



Sheffield Robotics Seminar

17 October 2023

slide 23



23



Talking to a Robot



Huang, G. & Moore, R. K. (submitted). Freedom Comes with a Cost?: How Affordance Designs Affect Users' Experience with a Conversational Social Robot. *Frontiers in Robotics and AI*.



Sheffield Robotics Seminar

17 October 2023

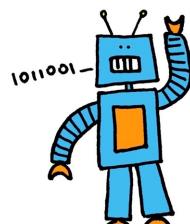
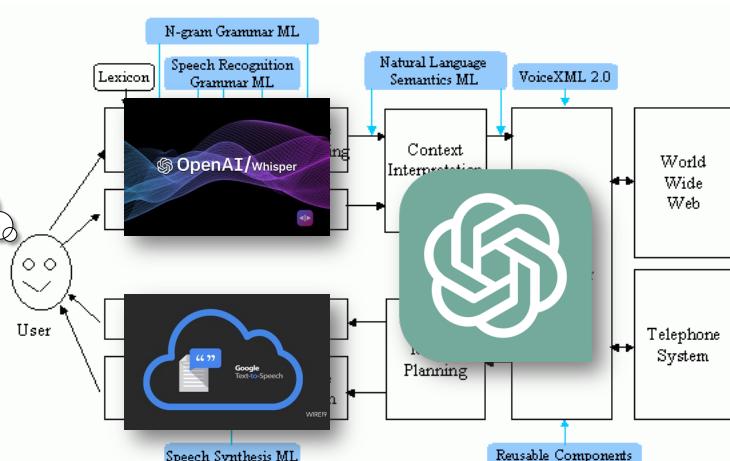
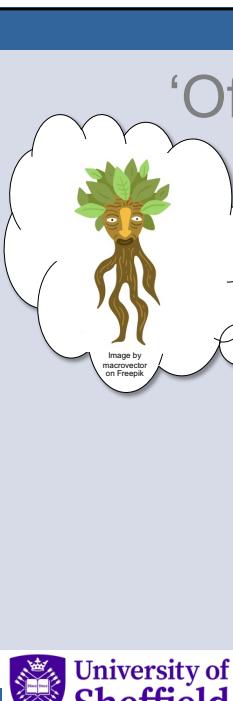
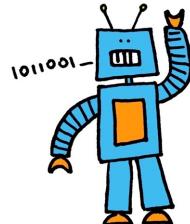
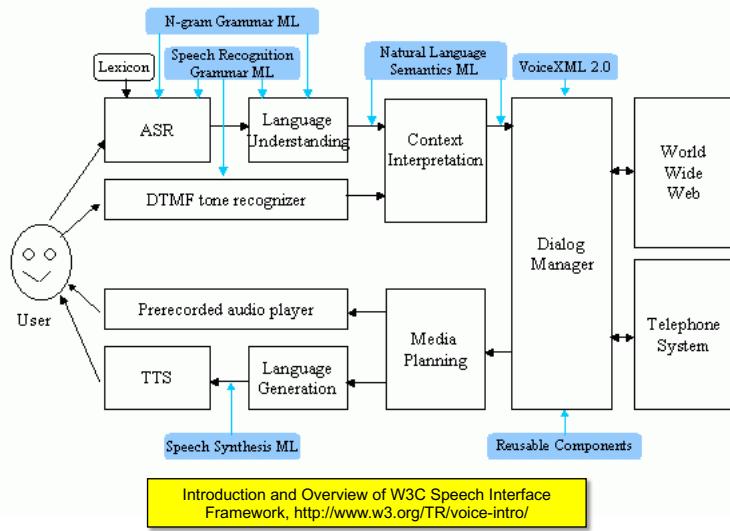
slide 24



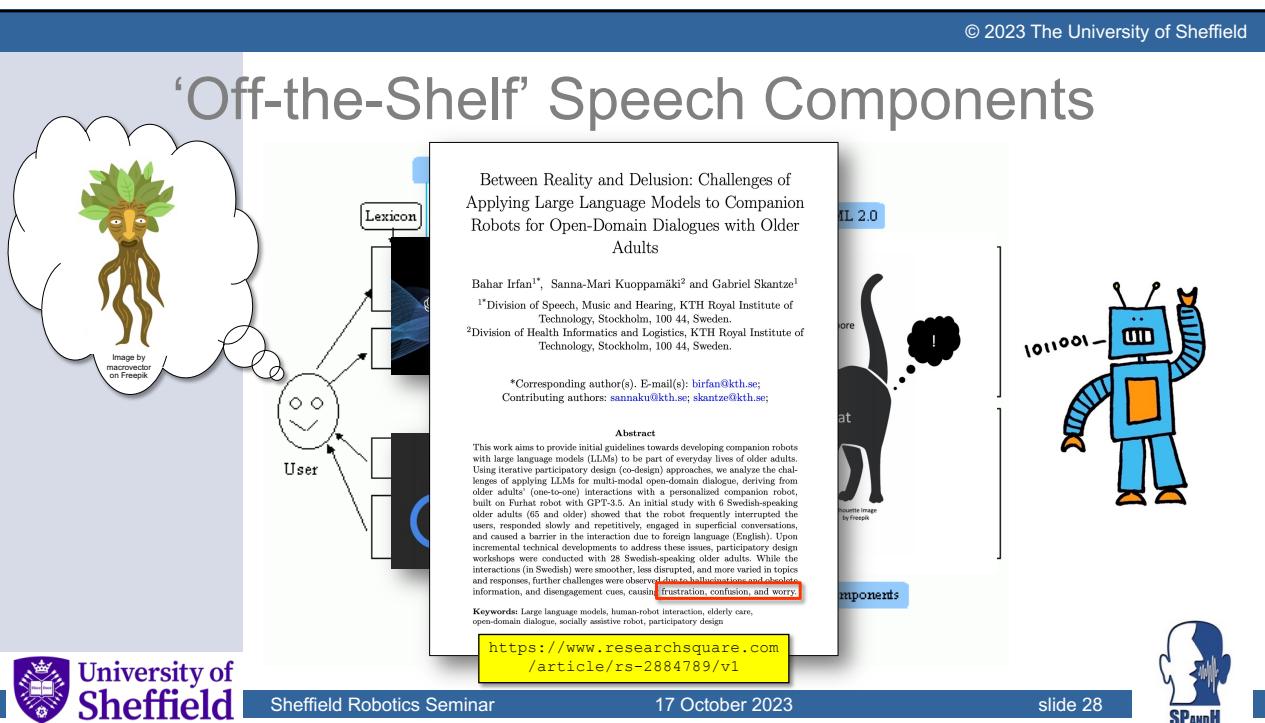
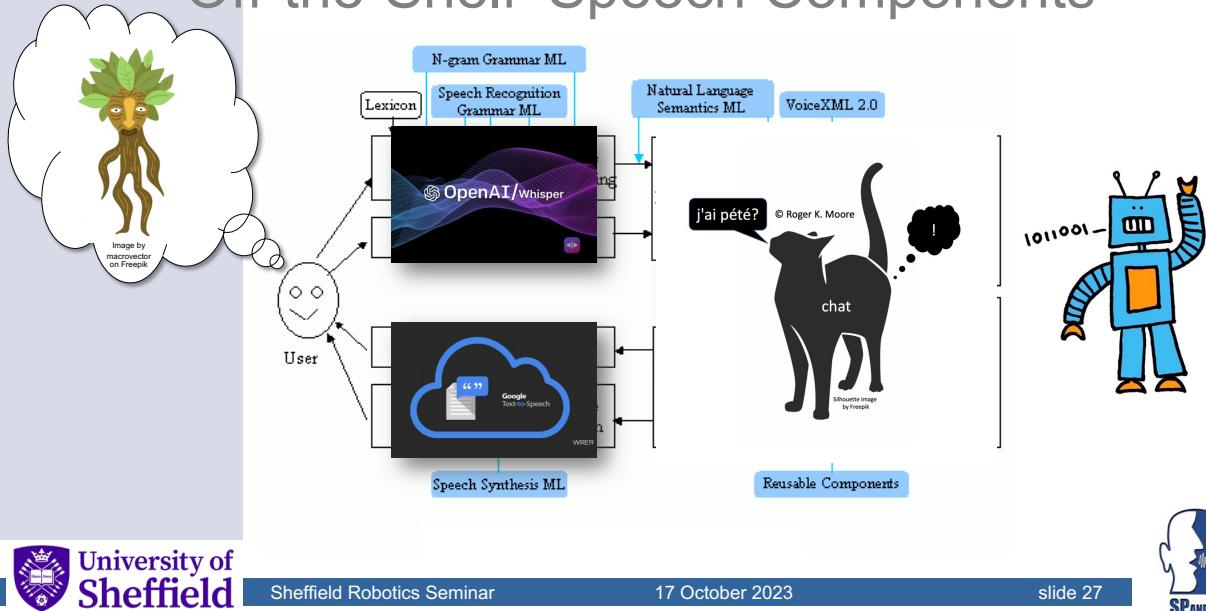
24

12

'Standard' Speech Interface



'Off-the-Shelf' Speech Components



Multi-Modal Communication

available



INFORMATION RATE

speech
hand gesture
pointing
eye gaze
facial expression
vocalisation
head movement
body pose
body movement

Moore, R. K. (2021). Embodied versus disembodied conversational agents In 5th International Workshop on Chatbot research: CONVERSATIONS-2021, Oslo.



Sheffield Robotics Seminar

17 October 2023

slide 29



29

Multi-Modal Human-Human Interaction

available



INFORMATION RATE



individuality
personality
demographics
culture

Moore, R. K. (2021). Embodied versus disembodied conversational agents In 5th International Workshop on Chatbot research: CONVERSATIONS-2021, Oslo.



Sheffield Robotics Seminar

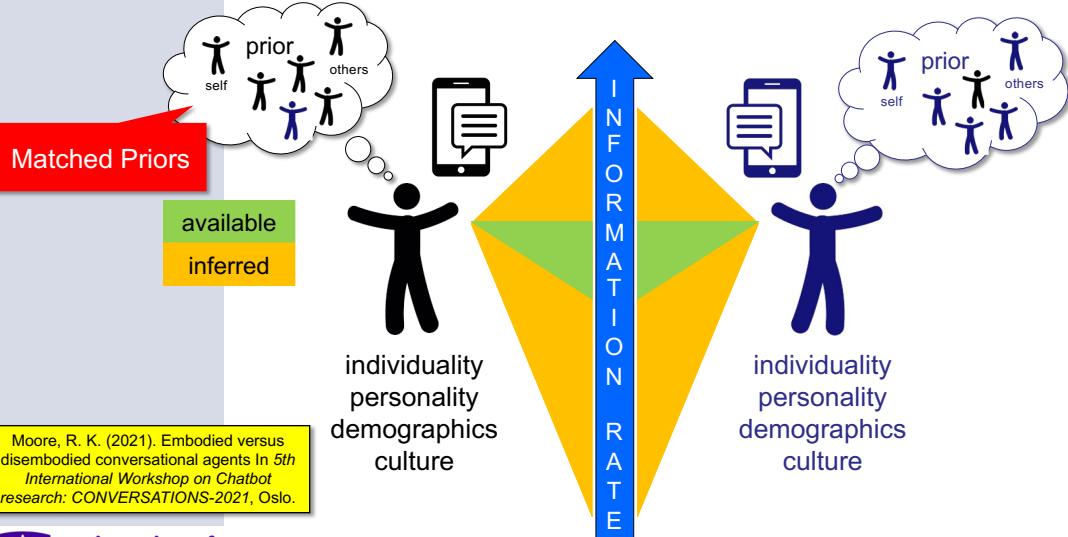
17 October 2023

slide 30



30

Mono-Modal Human-Human Interaction



Sheffield Robotics Seminar

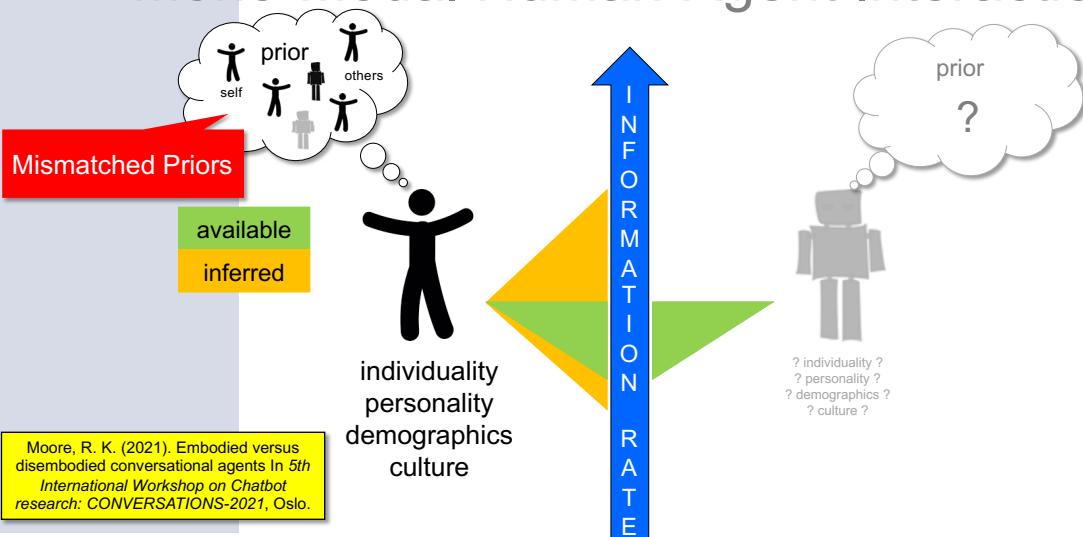
17 October 2023

slide 31



31

Mono-Modal Human-Agent Interaction



Sheffield Robotics Seminar

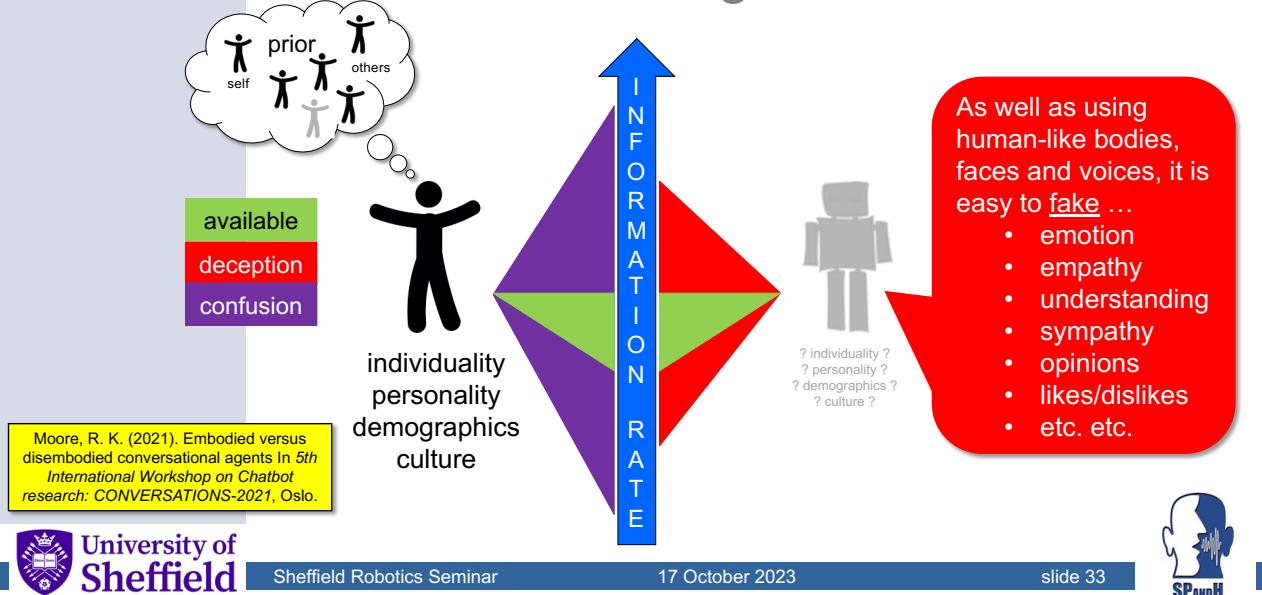
17 October 2023

slide 32



32

Multi-Modal Human-Agent Interaction



33

© 2023 The University of Sheffield

Mismatched Interlocutors

• Spoken language interaction between human beings is founded on **shared experiences, representations and priors** (*i.e. mutual situatedness/embodiment*)

• So, is there a fundamental limit to the language-based interaction that can take place between **mismatched interlocutors**?

Moore, R. K. (2016). Is spoken language all-or-nothing? Implications for future speech-based human-machine interaction. In K. Jokinen & G. Wilcock (Eds.), *Dialogues with Social Robots – Enablements, Analyses, and Evaluation*. Springer Lecture Notes in Electrical Engineering (LNEE).

University of Sheffield

Sheffield Robotics Seminar

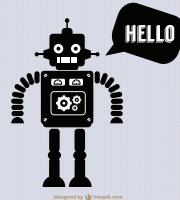
17 October 2023

slide 34

SPANDH

34

Talking to a Machine



Usability

Philips, M. (2006). Applications of spoken language technology and systems. In M. Gilbert & H. Ney (Eds.), *IEEE/ACL Workshop on Spoken Language Technology (SLT)*

Add NL/Dialog



Structured Dialog

Like a Human



'Habitability Gap'



Flexibility

17 October 2023

slide 35



Sheffield Robotics Seminar



35

SCIENTIFIC REPORTS

OPEN
A Bayesian explanation of the 'Uncanny Valley' effect and related psychological phenomena
Roger R. Moore

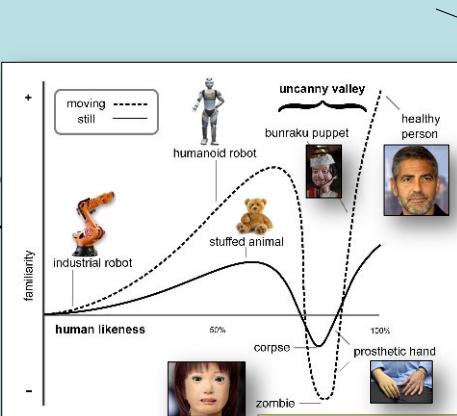
Received: 14 June 2012
Accepted: 9 October 2012
Published: 16 November 2012
DOI: 10.1038/srep02864
Cite this article: Moore, R. K. (2012). A Bayesian explanation of the "Uncanny Valley" effect and related psychological phenomena. *Nature Scientific Reports*, 2(864).

Moore, R. K. (2012). A Bayesian explanation of the "Uncanny Valley" effect and related psychological phenomena. *Nature Scientific Reports*, 2(864).

Structured Dialog

Add NL/Dialog

Like a Human

Mori, M. (1970). *Bukimi no tani (the uncanny valley)*. *Energy*, 7, 33-35.

Flexibility

17 October 2023

slide 36



Sheffield Robotics Seminar



36

© 2023 The University of Sheffield

The Uncanny Valley

Lead the conversation
Subscribe today and navigate your world with confidence

Opinion Virtual and Augmented Reality
What happens when AI passes through the 'uncanny valley'?
Robots are close to being so convincing that we can't tell them apart from humans — and that could be a problem

Why can't Nvidia boss Jensen Huang escape the Uncanny Valley that makes AI feel icky?
Is he human? Is he an avatar? Does it really even matter?

ARTnews
Gagosian's DALL-E-Enabled Art Exhibition
Throw Us Headfirst into the Uncanny Valley

IGN
Starfield character smiles are creepy and a game developer explains exactly why

NATIONAL GEOGRAPHIC
The uncanny valley, explained: Why you might find AI creepy
Even if we can't "see them", programs like ChatGPT's ability to emote like a human gives some people an eerie feeling. Scientists have a few theories on why this phenomenon happens.

IFLSCIENCE
The Uncanny Valley – What Is It?
The range of things that can produce this uncomfortable feeling seems to grow as our technology advances.

John Thorburn + Add to my

LOGIN Newsletters SUBSCRIBE

BENNETT MILLER United, 2022-23 Pigment print of AI-generated image BENNETT MILLER COURTESY THE ARTIST AND GAGOSIAN

50
25
Jan 1, 2004 Nov 1, 2009 Sep 1, 2015 Jul 1, 2021

Uncanniness ↓ ► affinity ↑ ► rapport ↑

University of Sheffield

Sheffield Robotics Seminar

17 October 2023

slide 37

37

© 2023 The University of Sheffield

A 'Canny' Approach

Moore, R. K. (2019). A 'Canny' Approach to Spoken Language Interfaces. *CHI-19 Workshop on Mapping Theoretical and Methodological Perspectives for Understanding Speech Interface Interactions*, Glasgow.

A 'Canny' Approach to Spoken Language Interfaces
Roger K. Moore
University of Sheffield
roger.moore@sheffield.ac.uk

ABSTRACT
Most embodied entities such as human-like are more familiar, but there appears to be a underlying point where they begin to trigger the 'uncanny valley' effect. This paper explores how agents can align the visual and vocal modalities to mitigate the negative effects of future well-established devices.

KEYWORDS
voice enabled devices, embodied, post uncanny valley effect, aligned affordances

INTRODUCTION
There has been a recent increasing interest in the development of voice enabled assistants. However, there is a lack of research investigating the potential benefits of these technologies. This paper aims to address this gap by investigating the potential benefits of aligning the visual and vocal modalities for voice enabled devices. The results show that aligning the visual and vocal modalities can reduce the negative effects of the 'uncanny valley' effect. This paper also highlights the importance of aligning the visual and vocal modalities for improving user experience and satisfaction with voice enabled devices.

CITATION
Moore, R. K. (2019). A 'Canny' Approach to Spoken Language Interfaces. In May (Ed.), *CHI-19 Workshop on Mapping Theoretical and Methodological Perspectives for Understanding Speech Interface Interactions*, Glasgow, UK.

uncanniness ↓ ► affinity ↑ ► rapport ↑

University of Sheffield

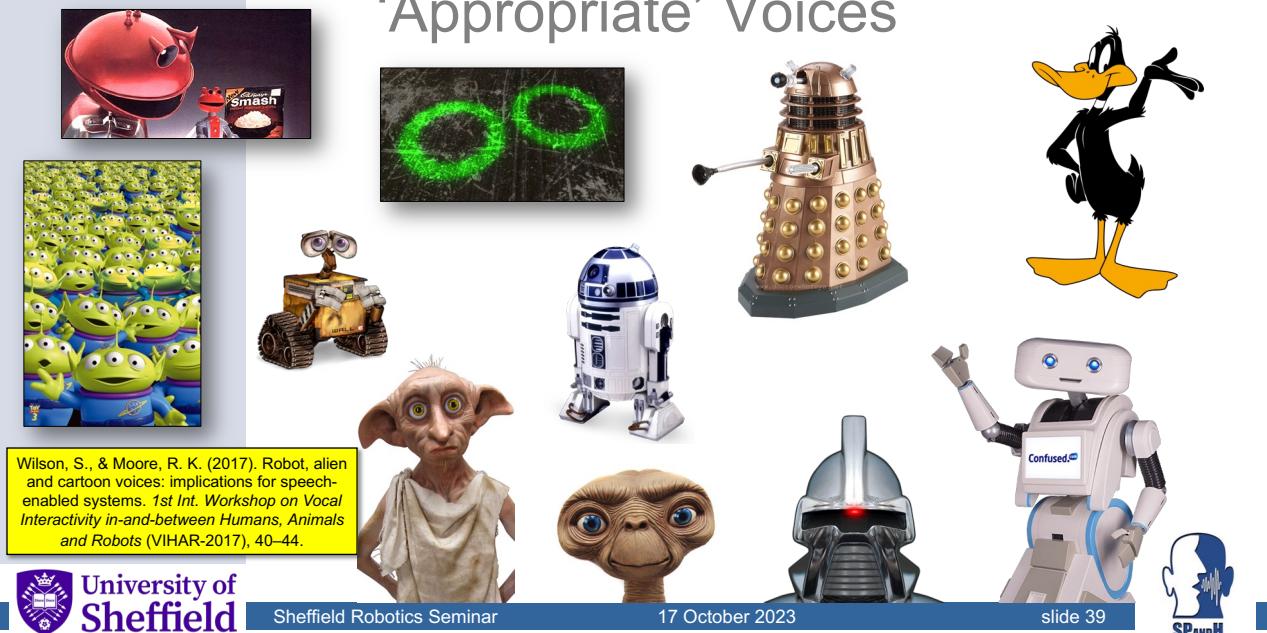
Sheffield Robotics Seminar

17 October 2023

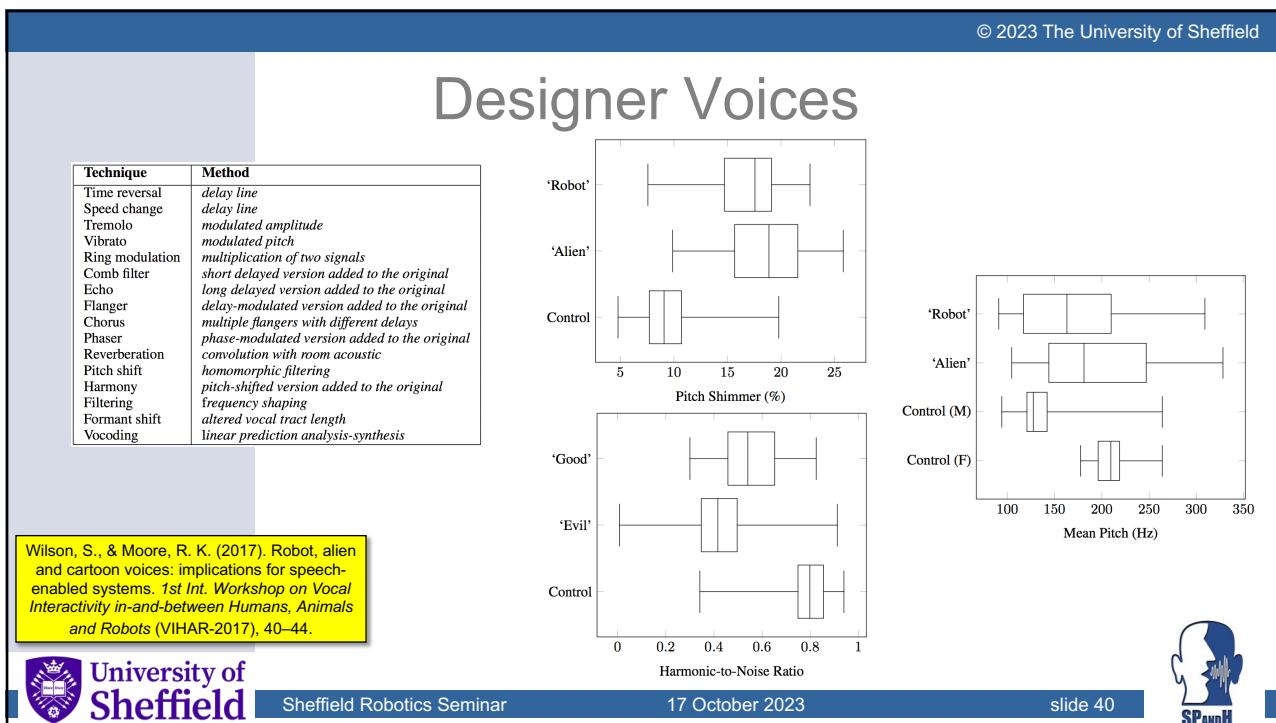
slide 38

38

'Appropriate' Voices



39



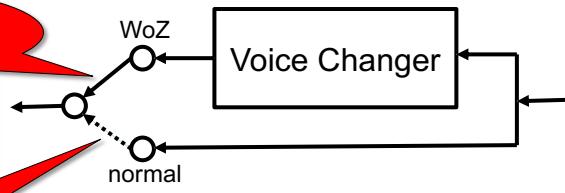
40

Deploying an ‘Appropriate’ Voice

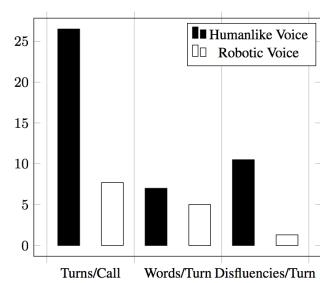
“Welcome to the route planning service - how can I help you?”



“Welcome to the route planning service - how can I help you?”



Moore, R. K., & Morris, A. (1992). Experiences collecting genuine spoken enquiries using WOZ techniques. 5th DARPA Workshop on Speech and Natural Language, 61–63.



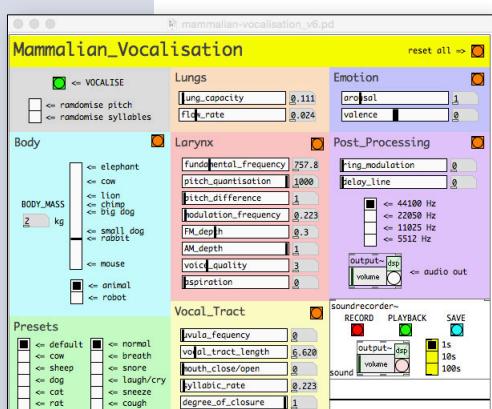
Sheffield Robotics Seminar

17 October 2023

slide 41

41

A Voice for the ‘MiRo’ Robot



Body mass = 2 kg
Breathing rate = 0.7 Hz
Fundamental frequency = 760 Hz
Vocal tract length = 6.6 cm

Moore, R. K., & Mitchinson, B. (2017). A biomimetic vocalisation system for MiRo. In *Living Machines 2017*. Stanford, CA.

Moore, R. K. (2016). A real-time parametric general-purpose mammalian vocal synthesiser. In *INTERSPEECH* (pp. 2636–2640). San Francisco, CA.



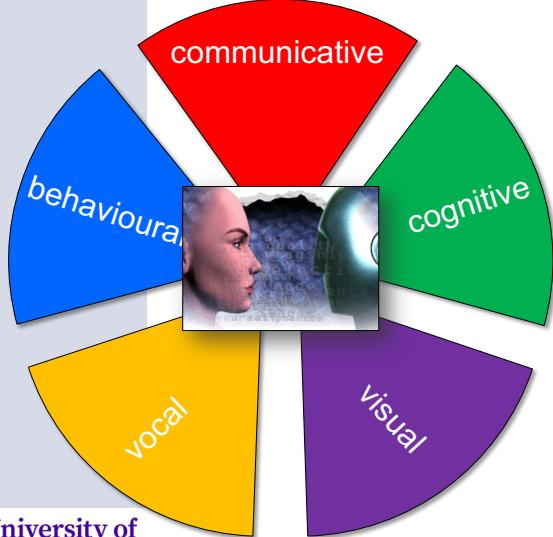
Sheffield Robotics Seminar

17 October 2023

slide 42

42

Bringing It All Together



Sheffield Robotics Seminar

Need to take a holistic approach to ...

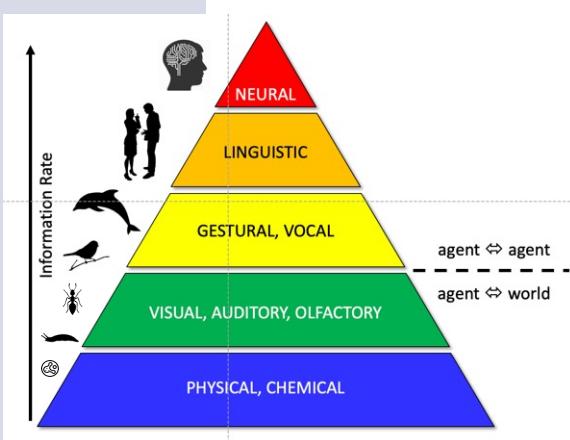
- designing and implementing interactional affordances
- accommodating user priors



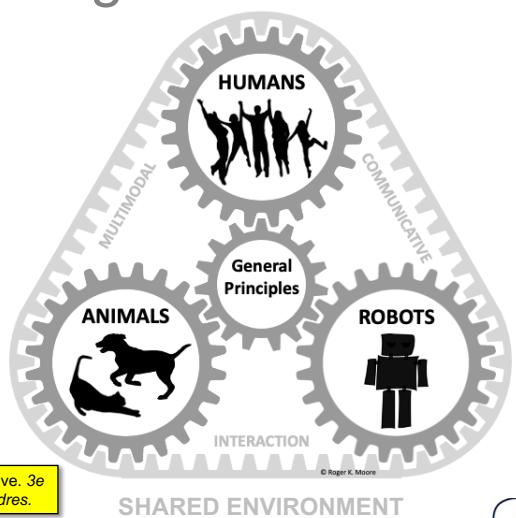
17 October 2023

slide 43

43



Moore, R. K. (2021). Où les êtres Sociaux: vers une théorie de l'interaction communicative. 3e Colloque International - Objets Animés, Humains, Animaux : Partenaires de Soins Tendres.



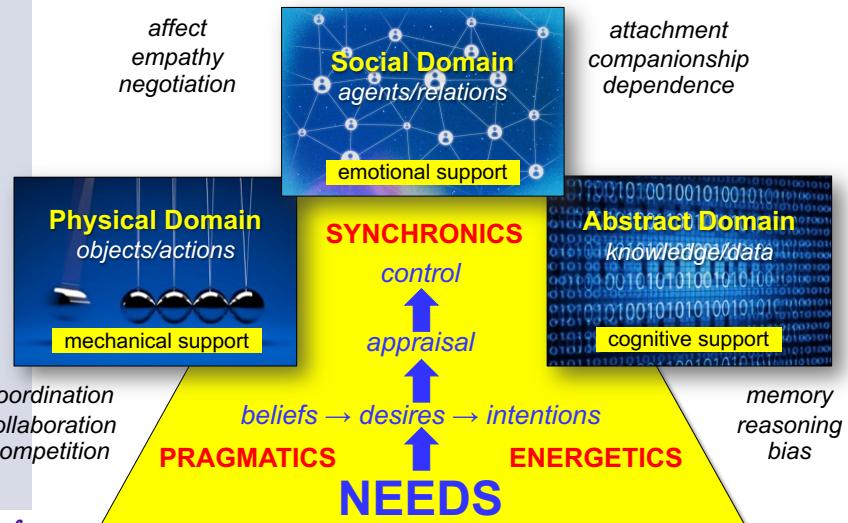
Sheffield Robotics Seminar

17 October 2023

slide 44

44

Bringing It All Together



Sheffield Robotics Seminar

17 October 2023

slide 45

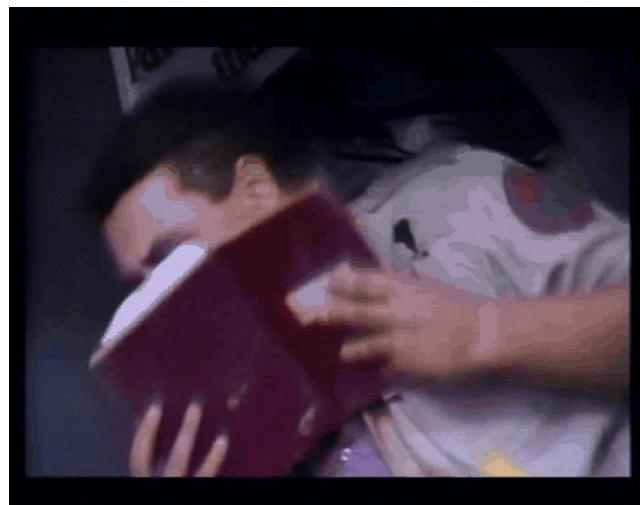


45

Bringing It All Together



<https://reddwarf.co.uk>



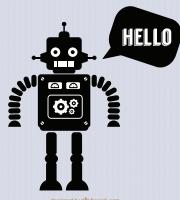
Sheffield Robotics Seminar

17 October 2023

slide 46

46

Summary



- Next-generation language-based interactive systems need to have a much **deeper understanding** of user behaviour
- Embodied systems need to facilitate **continuous adaptive multimodal incremental closed-loop coupling** between users and the physical, social and abstract worlds
- The visual, vocal and behavioural **affordances** of speech-enabled devices need to reflect their actual abilities
- We need to understand how language might function between **mismatched agents**
- We need to be able to measure (*and thus optimise and predict*) '**usability**' for language-enabled systems
- We must guard against **inappropriate use-cases** (e.g. infantilising users)



© 2023 The University of Sheffield

Where to Find Out More

<https://youtu.be/01ml3rpPZYg>

Sheffield Robotics Seminar

17 October 2023

slide 48

48

24



<http://staffwww.dcs.shef.ac.uk/people/R.K.Moore/>

49

© 2023 The University of Sheffield

Abstract

- It is often taken for granted that speech provides a ‘natural’ means by which users can interact with an autonomous agent such as a robot - especially if it is humanoid in design.
- Indeed, the ready availability of off-the-shelf spoken language tools such as Google’s speech-to-text/text-to-speech or OpenAI’s Whisper/ChatGPT makes it relatively easy to implement a ready-made solution.
- However, not only does such an approach risk creating a chimera - a potentially confusing and inappropriate compilation of misaligned visual, vocal and behavioural affordances - but it also facilitates a casual misrepresentation/faking of agency (e.g. by having the robot express apparently personal likes, dislikes and opinions using the pronoun “I”).
- This talk will address these issues, and provide an insight into why, when and how users might wish to converse with a robot.
- 60 mins = ~30 slides

 **University of
Sheffield** Sheffield Robotics Seminar 17 October 2023 slide 50 

50

25