

# D64- a Corpus of richly-recorded Conversational Interactions

Catharine Oertel, Fred Cummins,  
Nick Campbell, Jens Edlund, and Petra Wagner  
U. Bielefeld, University College Dublin, Trinity College Dublin, KTH, U. Bielefeld

*nick@tcd.ie*

# INTRODUCTION

- we recorded a corpus
- a VERY big corpus
- of VERY natural spoken interactions .....
- with lots of gear
- and now we are trying to cope with the mess!

# THE INITIAL MESS



# SPOKEN INTERACTION

- people always talk
- whenever they sit down together something social happens
- we are interested in how that happens (as a process)
- so we sat down and recorded ourselves
- with mocap (x6) audio (x12) video (x5) 360-degree (x2) etc
- and talked for 2 days!

# CORPUS COLLECTION



# PEOPLE & CONDITIONS

- there was 1 naive volunteer amongst 5 participants
- (ethics-committee-agreed release form were signed by all)
- NO constraints were imposed on the content of the conversations - but all participants have the right to ask for the removal of any section of the corpus at any time
- and we are crowd-sourcing the annotations . . . . .

# CHATTING OVER DRINKS



# A BRIEF HISTORY OF THE PROJECT

- atr/jst-crest esp corpus (Expressive Speech Processing)
- scope robot's ears project (round-table meetings)
- nov07 - friends sitting around a table in the lab
- aug09 - friends coming round for a sushi meal at home
- dec09 - friends in a room in dublin (drinkin' coffee & wine)
- fastnet! (focus on actions in social talk (network-enabling))



# THE RECORDINGS

- D64 is the hotel room where Jens was staying (in Dublin)
- we wired it up - it was a mess - then we hid the wires
- we wore radio mics and were free to move around
- there was a kitchen area so we also ate & drank ..... wine too!
- we recorded a morning and an afternoon session on day 1
- and an all-day session on day 2 ending with a meal (& wine)

# ALL-SEEING EYE



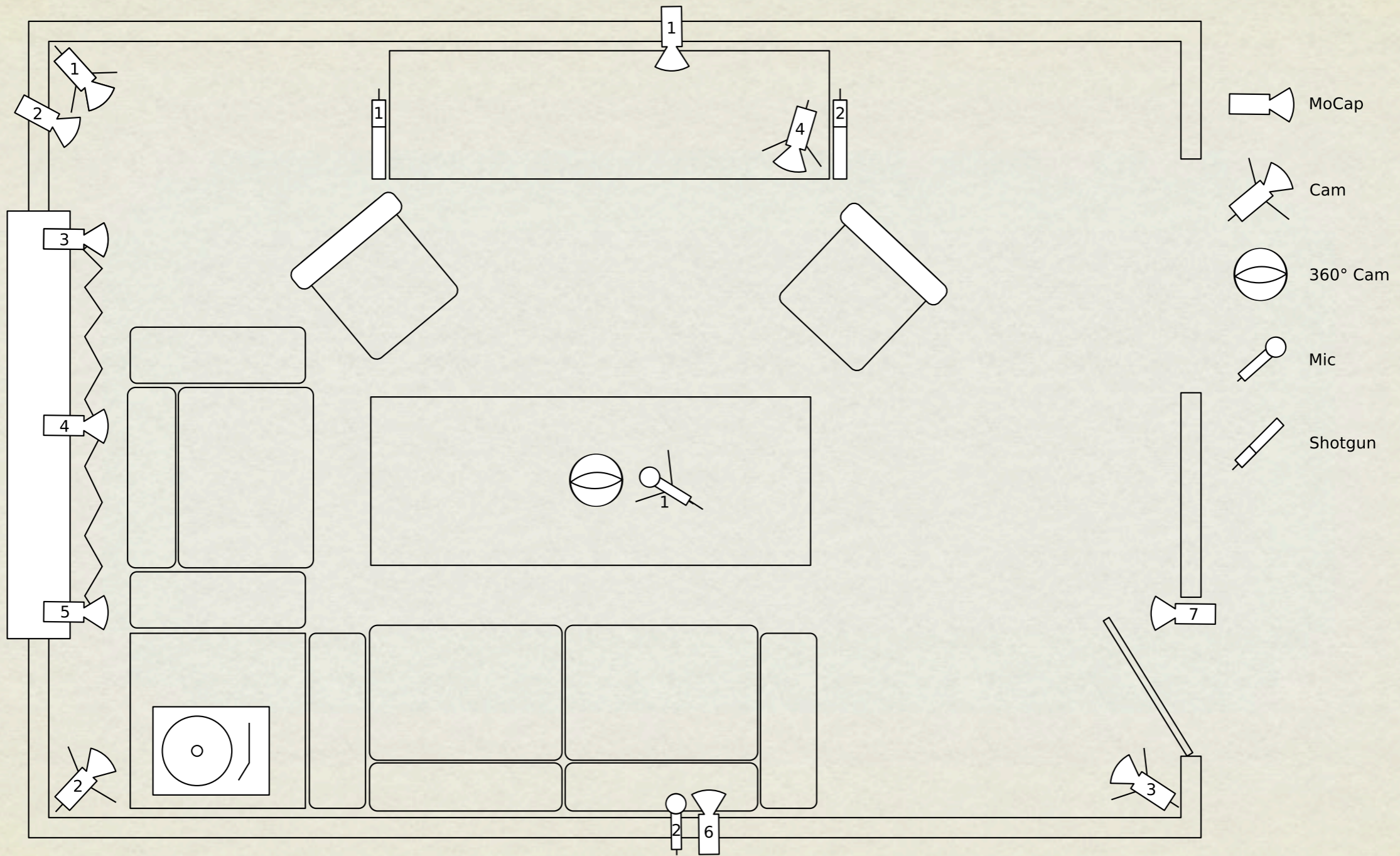
# WIRED UP! (TRACKED)



# CLEARED UP & WORKING



# THE ROOM LAYOUT



# ON A SCALE OF 'REALITY'



unconstrained  
covert field  
recordings

D64 corpus

spontaneous  
speech, controlled  
laboratory setting,  
e.g. "map tasks"

read speech,  
laboratory setting

- this is NOT read-speech nor scripted or prompted
- it is a spontaneous social interaction and has HIGH 'reality'

# SOCIAL INTERACTION



# TALKING WHILE EATING





# DATA HANDLING

## Multimodal Expressive data (d64)

### ◆◆◆audio◆◆◆◆

unless you really need the original recording, try the mp3 first - for easier downloads!

	day1 am		day1 pm		day2	
left field	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
right field	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
mix1	-	<a href="#">a b</a>	-	<a href="#">x</a>	-	<a href="#">a b c d</a>
mix2	-	<a href="#">a b</a>	-	<a href="#">y</a>	-	<a href="#">a b c d</a>
Nike	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	-	-
Catha	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
Jens	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
Fred	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
Nick	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>	<a href="#">(big)</a>	<a href="#">mp3</a>
table	???	???	???	???	<a href="#">(big)</a>	<a href="#">mp3</a>
sofa	-	<a href="#">a b</a>	-	<a href="#">y</a>	-	<a href="#">a b c d</a>
sony-b66	-	<a href="#">a b</a>	-	<a href="#">y</a>	-	<a href="#">a b c d</a>
head-mount					<a href="#">(big)</a>	<a href="#">nick</a>
sync	<a href="#">(big)</a>	-	<a href="#">(big)</a>	-	<a href="#">(big)</a>	-

these tables were filled from logic rather than listening ...  
from wiring diagrams rather than ears-on experience ....  
they are bound to be wrong!!

### ◆◆◆video◆◆◆◆

most of these are compressed to mp4 - should download MUCH faster than the originals

	day1 am		day1 pm		day2	
sofa	<a href="#">all</a>	<a href="#">cut</a>	-	<a href="#">cut</a>	-	<a href="#">a??</a>
armchair	<a href="#">a b c</a>	-	<a href="#">a b c</a> <a href="#">d e</a>	-	-	<a href="#">a b c d e</a> <a href="#">f g h i j k</a>
chair	<a href="#">uncut</a>	-	-	-	-	<a href="#">clearing up</a>
sofa left	<a href="#">all</a>	-	-	-	<a href="#">uncut</a>	-
fred	-	<a href="#">1 2</a>	-	<a href="#">3 4</a>	-	<a href="#">5 6 7 8</a>

see note above .. same applies here!

### ◆◆◆extras◆◆◆◆

	day1 am	day1 pm	day2
loops		<a href="#">(fred)</a>	
360-degree camera	<a href="#">a b c d e f</a>	<a href="#">a b c d e f g h i j k</a> <a href="#">l m n o p q r s t u v</a> <a href="#">w x y z aa ab ac ad ae</a>	<a href="#">a b c d e f g</a>
other round	<a href="#">1 2</a>	<a href="#">3a 3b</a>	<a href="#">4 5</a>
mocap	<a href="#">set 1</a>	<a href="#">set 2</a>	<a href="#">set 3</a>
4-track bkp audio	<a href="#">?1 ?2</a>	<a href="#">x1 x2 x3 x4 x5</a> <a href="#">x6 x7 x8 x9</a>	<a href="#">xa xb xc xd xe</a>

### ◆◆◆photos◆◆◆◆

see [here](#) (see separate page)

# AN ANALYSIS

- we were particularly interested here in ‘virtual spaces’
- we were all occupying the same physical space, but in different cognitive spaces - or attentional states
  - can distances be measured in these spaces?
  - does ‘proximity’ correlate with ‘arousal’?

# AROUSAL & DISTANCE

- we hope to measure the ‘push&pull’ of social interaction . . .  
(cf Murray & Trevarthen - mothers & babies over video)
- what are the dynamics of spoken interaction?
- are there fundamental ‘laws’ or ‘forces’ that drive it?
- AROUSAL: when the group gets ‘heated’ - what happens?
- SOCIAL DISTANCE: when two people interact - is it ‘focus’?

# THE EBB & FLOW OF JOINT INTERACTION

- simultaneous & real-time movements of several participants at once - like puppets being pulled by a common string . . .



# SUMMARY

- this is a massively multi-modal 8-hour corpus
- this large-scale over-sampling provides very rich data
- it is available for research purposes (contact [nick@tcd.ie](mailto:nick@tcd.ie))
  
- web-based interfaces are being developed and we hope that the availability of this multi-faceted data will encourage further research into the mechanisms of spoken-interaction

# CONCLUSION

- this paper has presented a new corpus of multimodal data
- it features spontaneous multi-party social interaction
- captured in a multi-featured way - motion, audio, and video
- we are currently assembling and annotating the material
- it will be made available under [www.speech-data.jp/mmx](http://www.speech-data.jp/mmx)
- MMX is Multi-Modal-Expressive (2010 in roman numerals)

# ACKNOWLEDGMENTS

- This work has been supported by grants to Nick Campbell from the Visiting Professorships & Fellowships Benefaction Fund from Trinity College Dublin, and the Kaken-B Fund for Advanced Research from the Japanese Ministry of Information, Science & Technology, and also Science Foundation Ireland, Stokes Professorship Award 07/SK/I1218. Jens Edlund is supported by The Swedish Research Council KFI Grant for large databases (VR 2006-7482). Catharine Oertel is supported by the German BMBF female professors programme (Professorinnenprogramm) awarded to Petra Wagner.
- Finally, thanks to Nike Stam for her generous participation ;- )

THANK YOU  
FOR LISTENING