

Transcript: “Smash the vase”

From Anthropic (YouTube Channel),
episode “What should an AI's personality be?” (08.06.2024),
36:21 bis Ende
URL: <https://www.youtube.com/watch?v=iyJj9RxSsBY>

- Amanda Askell (Alignment Finetuning Researcher at Anthropic)
- Stuart J. Ritchie (Research Communications at Anthropic)

This is the end of an expert interview produced by the US-based AI startup company [Anthropic](#) for promotional purposes. Anthropic prides itself to be a company doing “AI research and products that put safety at the frontier”¹. Its main offering is the Large Language Model (LLM) “Claude”.

Amanda Askell is a well-known expert for AI ethics and AI alignment and serves in various interviews (podcasts and YouTube) as the public figurehead of Anthropic for these topics. Stuart J. Ritchie is a Scottish psychologist also working for Anthropic doing “research focused comms”².

The recording of the interview is 37:40 minutes long. Mr Ritchie clearly takes the role of explaining and supplementing Ms Askell’s statements for non-expert listeners. He does this in a rather forceful manner, often strongly interrupting Ms Askell. His backchannels are particularly loud and hearably disturb Ms Askell’s talk at some points. So much so that one YouTube user comments “Hey, Stuart. I'd like to hear Amanda finish a thought, if you don't mind.” With 4 other users chiming in with the same sentiment.

Research questions

I am interested in extended trouble-induced side sequences in interviews launched by the interviewer dealing with problematic (offensive, delicate, mispronounced etc.) utterances. In this excerpt, I am especially interested in:

- How does the sudden topic change during an expert interview play out:
 - Orientation to and handling of the intervention
 - Shift in epistemic stance and status of the two interactants
 - Attempts to return to the institutional context and the current interview topic
- How is the “issue” acknowledged and the “wrong-doing” corrected?
- Is there a recognizable release of tension during and at the end of the intervention?
- The role of overlap and competing for the floor
- Intonation and laughter marking the handling of trouble and changes in stance
- Gender, Identity, Personality: How do they come into play in this excerpt?

¹ See homepage of anthropic.com (accessed 28/03/2025)

² Twitter: <https://x.com/StuartJRitchie/status/1744764822011699553>

01 Ama >If you were to< sho: ↑excessive empathy .h
02 >eemee ya could imagine like someone shewin excessive empathy< to
03 leyk: (.) to objects in the #wɔ:rld#. .h >and being [like,]<
04 Stu [>Ya<]
05 Ama Q:h >you-<you should go to pri:son: if you like if you smash the
06 vayz ↑>an I'm like< loo:k I think it's↑ good to not get in the habit
07 of like [h smashing objects.]
08 Stu [S:orry can I just stop] you [there]?
09 Ama [>ya ya<]
10 Stu You:re Sco:ttish and you just said vay:z.
11 (.)
12 Stu That'[s-]
13 Ama [Q:h] is [that?]
14 Stu [er:]::: >You kent say that< [it's vah]z
15 Ama [iszat]
16 =>Iszat Amer[ican?<]
17 Stu >[Ye-high long] hv you been in America?<
18 Ama £[I've been in America for like thirteen] yearsh£ hh h h h h
19 Stu [Yea, >How long have you been in America?]
20 Stu T00: LQ:NG [clearly it's been too long]
21 Ama [.hhi ha ha ha]
22 Stu =B'cause you say vay:z
23 Ama =↑£Smash the vayz on the sidewalk£↑ HH [h h]
24 Stu [>Carry on,<]
25 Stu ↑Say-ɔ:h my ↓God ɔ:h
26 (.)
27 Ama .hh [£ɔ:h->I neva-even<↑]
28 Stu [(>terrible< I k-)]
29 Ama =>Uh-uh<-Ive-f:orgotten things £that are Scottish and aren't, so£
30 [didn va:z, okay, va:z heh]
31 Stu [Hyea well, we cert'nly don't] say-
32 =no one in this country's ever said vay:z. >Let me [tell ya.< (.)]
33 Ama [↑£Okayokay£↑]
34 Stu A:h
35 Ama =smash [the va:z.]
36 Stu [>↑Pleas-carry-on!<]
37 (.)
38 Stu [>Pleas-carry-on,<]
39 Ama °[yea-okay]
40 =s:o-like u::hm°(.)
41 Yea, ↑so >i-if ya wa to say to peo:ple<↑ o:h you should like go to
42 pri:s'n: fer fer smashing a vaz £u::hm: hh the:n £
43 Stu =>sankya<
43 Ama like ghh h h .hhih tha:t >that's gone too far<, =So there is like
44 risks on all sides here,
45 =But ↑yeah >[maybe I am] sympathetic to the idea of like< (.)
46 Stu [<yess>]
47 Ama Do:nt-n like nee:dlessly like lie to or-or mistreat

48 Stu >°right°<
49 Ama =treat (.) <anything and that kind of includes these things
50 >°even if you think they're #not moral patients#°<.