# Language in Interaction
## Speaker and Listener information

SS16 - (Embodied) Language Comprehension

Maria Staudte

# So far …

- Embodiment

- Situated & embodied language learning

- Situated adult language comprehension (& production)

- Language in Interaction

  - Taking another person into account

  - Sending and perceiving bodily signals

# Language in/for Interaction

- Presupposes a listening/speaking partner

- Both partners use more than spoken language

  - Non-verbal signals: Facial expression, emotions, gaze, posture, gesture etc

- How do these influence language processes?
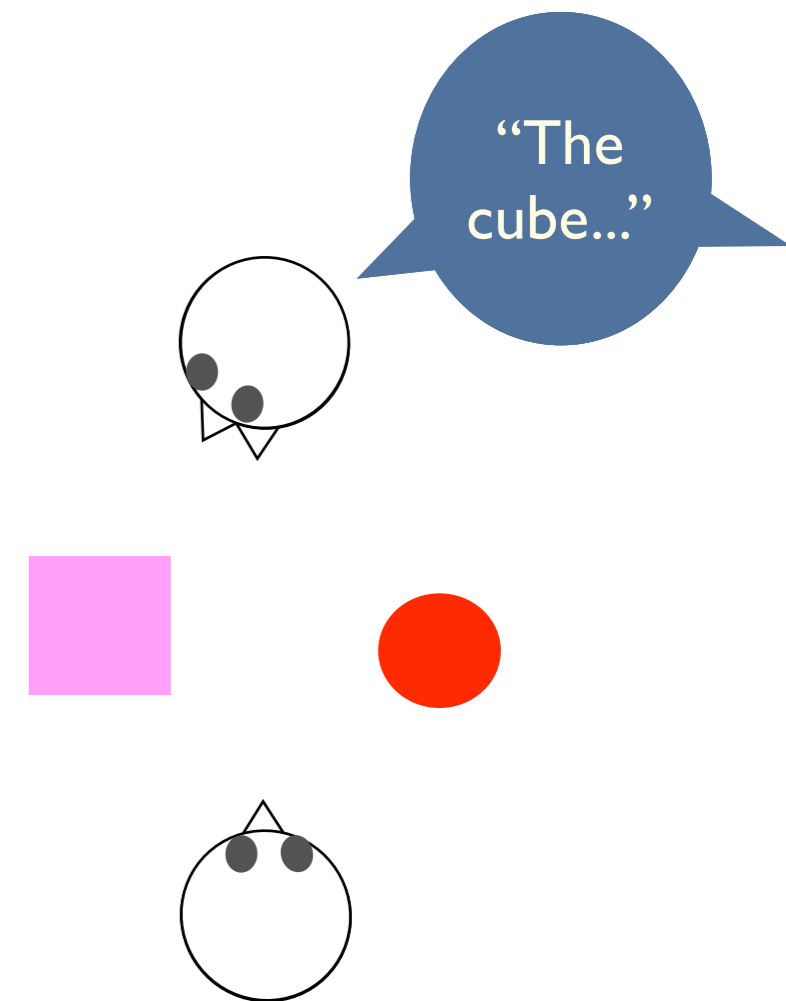
  - Information contribution, timing, cost

# Gaze

- Eye-movements reflect comprehension/prediction/planning processes

  - Measure

- Eye-movements are a signal by themselves to the partner!

  - Speaker gaze

  - Listener gaze

# Speaker gaze

# Referential Gaze in Communication

- Speakers look at what they are about to mention (e.g. Griffin & Bock 2000)

- Listeners look at what they hear (e.g. Tanenhaus et al., 1995)

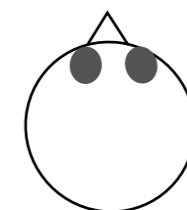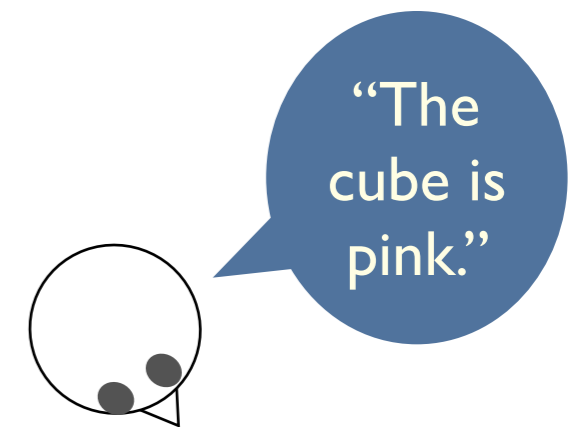- Listeners look at what the speaker looks at (e.g. Hanna & Brennan 2007)
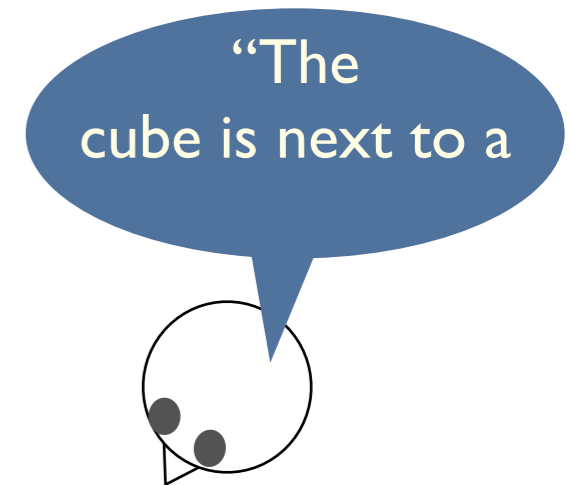
Jim Botsacos in "Cooking together"

# Speaker Gaze

✦ Listeners follow speaker gaze & utterance

✦ Facilitation/Disruption effect on sentence validation (congruent vs incongruent)

✦ Temporal shifts are irrelevant

✦ Cause of these effects?

"The cube is pink."

# Visual Attention & Order

- Speaker gaze & utterance both provide cues that drive listeners' visual attention

- Is order relevant for the utility of cues?

- Manipulate cue order to:

  - Explore information integration

  - Shed light on the role of information provided by speaker gaze

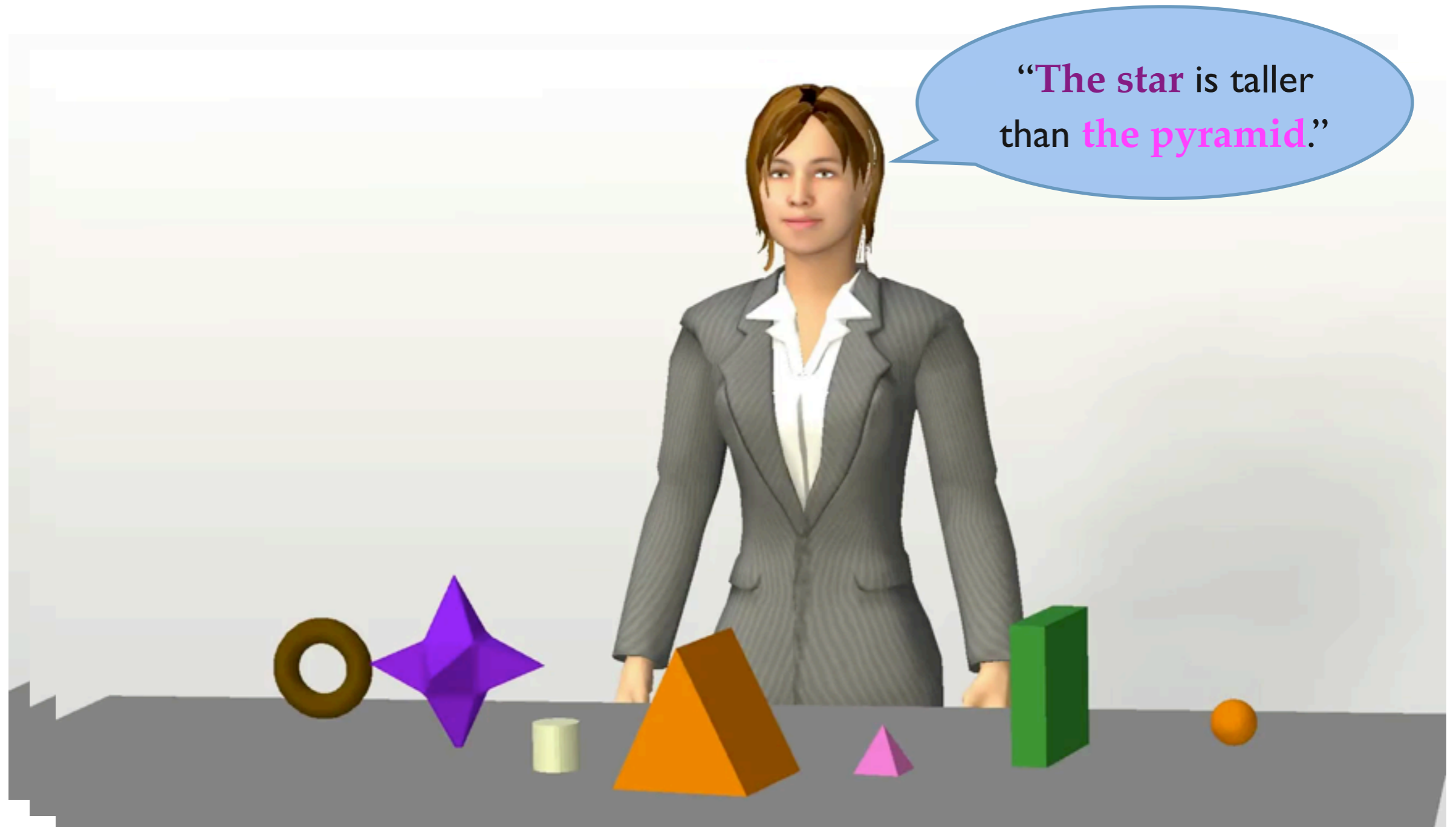# Visual Attention & Order

- Speaker gaze & utterance both provide cues that drive listeners' visual attention

- Is order relevant for the utility of cues?

- Manipulate cue order to:

  - Is there a bias towards preferring one modality?

  - Shed light on the role of information provided by speaker gaze

# Visual Attention & Order

✦ Speaker gaze & utterance both provide cues that drive listeners' visual attention

✦ Is order relevant for the utility of cues?

✦ Manipulate cue order to:

> ✦ Is there a bias towards preferring one modality?
>
> ✦ Is gaze like any other visual cue, simply increasing visual saliency? (persistent highlighting)
>
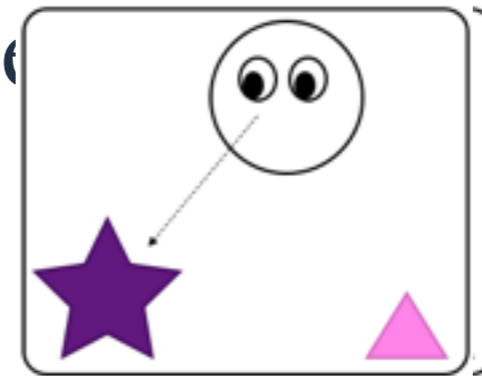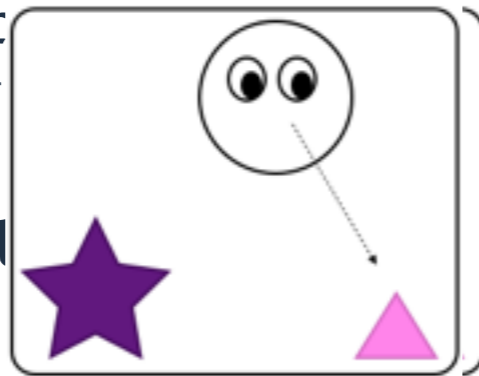> ✦ Does gaze signal intentions linked to utterance?

# Experiment 1

✦ Task: Is the utterance correct or not?

✦ 3 Conditions:

  ✦ Congruent, Reverse, Neutral

✦ 36 diff[...]ts (12 shape[...]

✦ All cou[...]ed

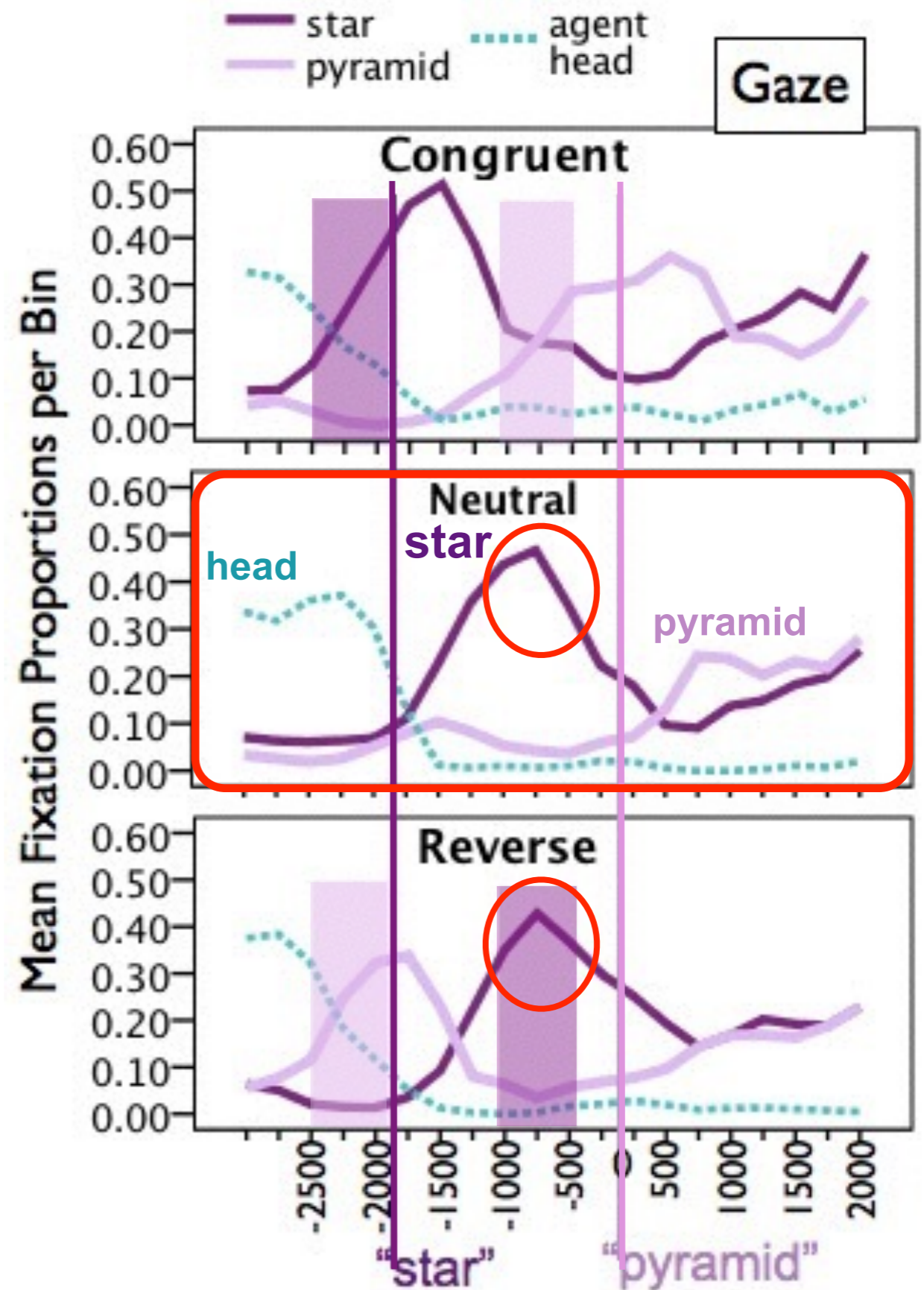✦ DV: Eye-movements, Response time

Reverse: Congruent: *The **star** is taller than the **pyramid***

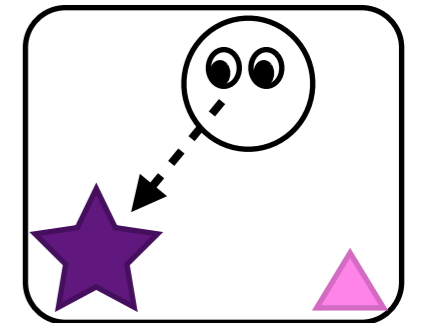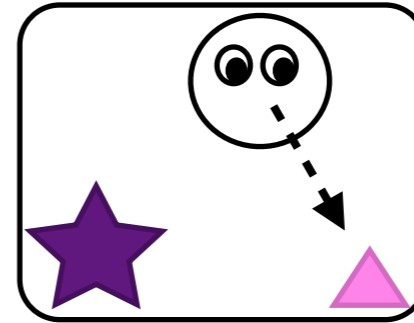# Experiment 1 : Results

# Experiment 1 : Results

- Eye-movement data:

    - Visual attention shifts are elicited by both speaker gaze and utterance, possibly automatically

- Response time data:

    - Visual information, gained through speech- and gaze-mediated attention shifts, is integrated

    - Information integration is difficult in reverse condition!

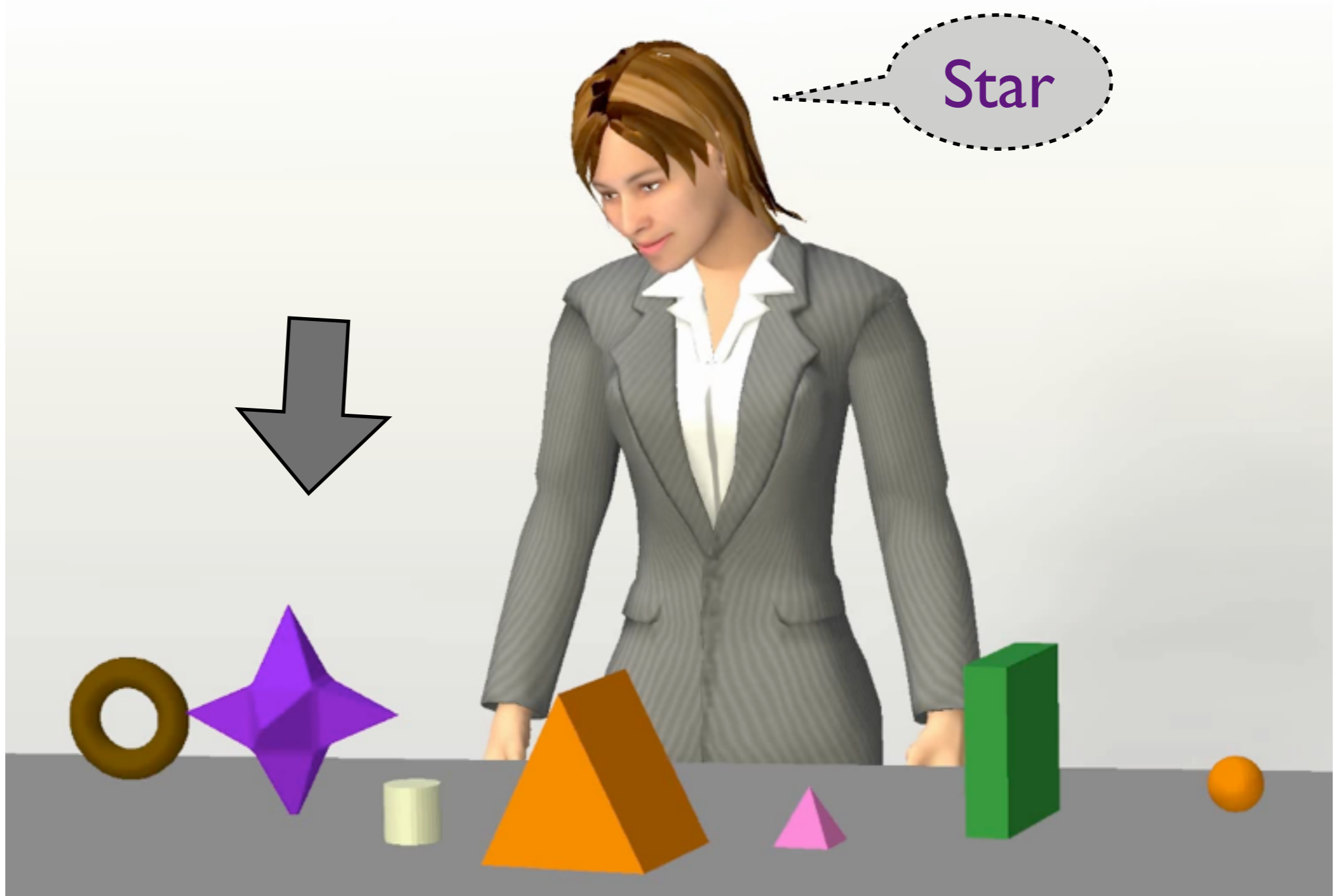# Experiment 1 : Results



*"The star is taller than the pyramid."*
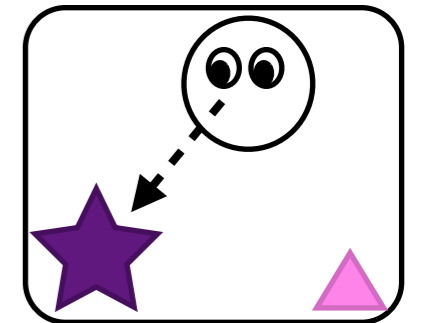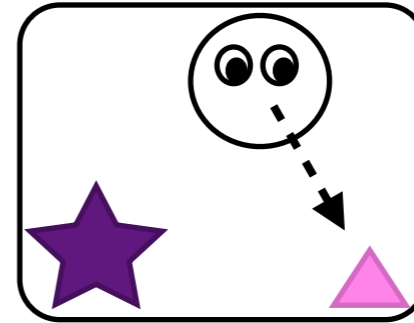
✦ Reverse cues:

  ✦ Gaze-mediated fixations to "pyramid"

  ✦ Speech-mediated fixations to "star"

➡ **RT data reveals disruption instead of facilitation!**

✦ What causes the slowed response time?

# Experiment 1 : Results



*"The star is taller than the pyramid."*

- Reverse cues:
  - Gaze-mediated fixations to "pyramid"
  - Speech-mediated fixations to "star"

  ➡ **RT data reveals disruption instead of facilitation!**

- What causes the slowed response time?
  - Timing and saliency? Or referential intentions?

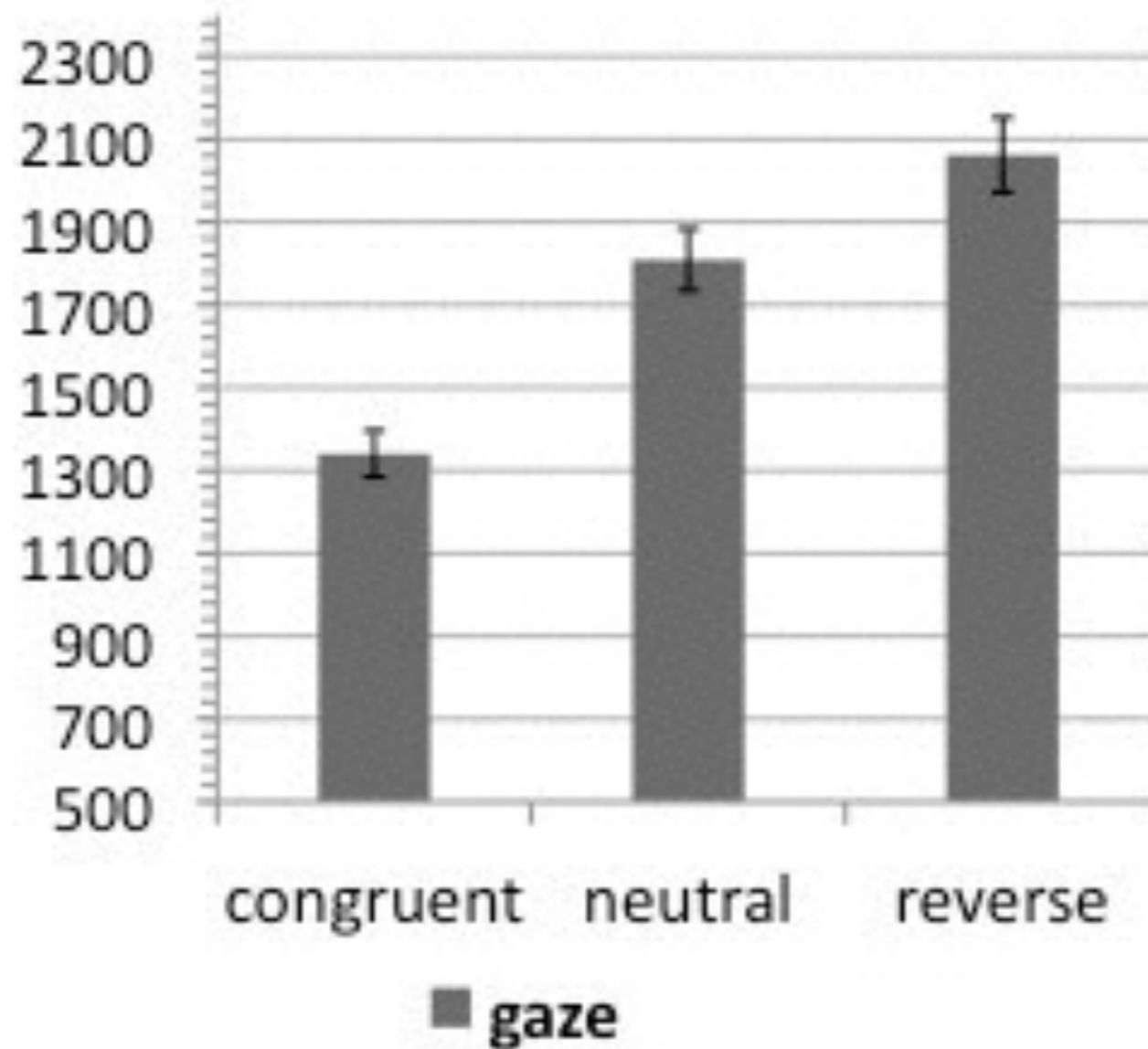- Do other (purely) visual cues have the same effect?

# Experiment 2

# Experiment 1: Results



**Response Times in 3 Conditions**

gaze

# Experiment 1+2 : Results



**Response Times in 3 Conditions for each Cue Type**

# Experiment 1+2 : Results

# Experiment 1+2 : Results



**Response Times in 3 Conditions for each Cue Type**

Cue Type x Condition

# Experiment 1+2: Results

# Further Results



✦ Response Time Block Analysis:

  ✦ Learning effect for <u>reverse</u> arrows (interaction)

  ✦ No learning effect for <u>reverse</u> gaze (no interaction)

# Interim summary

✦ Gaze elicits a prediction for the next referent

✦ Strong bias to infer referential intentions as acquired across many years

✦ Arrows are assigned a task-specific utility

✦ Unbiased cue which can be used flexibly

✦ Gaze affects comprehension *beyond* visual cueing

# Remaining issues

- Difference between speaker gaze and arrow cues:

  - Precision of cue

  - Reliability wrt language

# Precision?

- Gaze is (often?) less precise than e.g. arrow cues

- Compare arrows against simplified, precise gaze cue

  ➡ Benefit in "reverse" condition!



Response Times in 3 Conditions for each Cue Type

■ gaze  ■ arrow



Response Times in 3 Conditions for Arrows and Precise-Gaze

■ arrow  ■ precGaze

# Reliability

- Gaze occurs more often / more naturally than arrows

- Tendency to trust & follow gaze more than arrows?

  - Arrow usage more strategic?

  - Compare 0% (as before) with 25%, 50%, 75% trials with **invalid** cues in experiment

# Reliability - Cue following



Congr.

"star"  "pyramid"

No

Rev.

0%          25%          50%          75%

(Embodied) Language Comprehension                    Speaker/Listener

# Reliability - Cue effect

# Reliability

- Listeners stop following gaze (only) when cue is misleading in 75% of trials

- Listeners keep following and benefiting from arrows (even when cue misleads in 75%) !

➡ Gaze-following is less automatic

➡ But also less strategic than arrow usage

# Listener gaze

# Listener Gaze

- Listeners look at

  - what they hear

  - what the speaker looks at

- Speakers monitor what listeners look at

- How can/do they exploit this information?

  - Can we evaluate instructions (better) using eye-tracking?

  - Can we construct instructions (better) using eye-tracking?

# The Task



GIVE setting (Koller et al., 2010)

# Recording object inspections



**faceLAB eye-tracking system**

- Every 15ms, sample 2D position on screen that the user is fixating

- Resolve this position to the corresponding object in the current 3D scene

- User looks to an object of more than 300ms count as a "referential inspection" of that object

# Tracking listener gaze

# Monitoring understanding

- Based on eye gaze, system attempts to predict whether the user has understood its referring expressions

- System generates proactive feedback accordingly

  - Target inspection: "Yes, that one!"

  - Distractor inspection: "No, not that one!"

# Setup

# Example scene

"Push the left button to the..."

# Example scene

"...right of the flower.      "

# Example scene

"...flower.     - Yes, that one."

# Baseline 1:
# No feedback

- No monitoring of referential understanding

- No proactive feedback

- System generates a follow-up referring expression only after user has pressed wrong button or asked for help

# Baseline 2:
## Movement-based feedback

- System makes prediction only if user moves towards single visible button

- Same feedback as gaze-based system

  - Movement towards target: "Yes, that one!"

  - Movement towards distractor: "No, not that one!"

# Eye movements

|  | inspection durations | |
|---|---|---|
|  | target | distractor |
| **successful** | | |
| eyetracking | 2111.6 | 720.5 |
| no-feedback | 1492.0*** | 185.7*** |
| movement | 1493.8** | 260.5*** |
| **unsuccessful** | | |
| eyetracking | 752.1 | 3378.9 |
| no-feedback | 619.5 | 1891.7 |
| movement | 602.6 | 2113.1 |

*Differences to eyetracking system statistically significant at ***$p < 0.001$, **$p < 0.01$, *$p < 0.05$*

# Eye movements

|  | inspection durations | |
|---|---|---|
|  | target | distractor |
| **successful** | | |
| eyetracking | **2111.6** | 720.5 |
| no-feedback | **1492.0\*\*\*** | 185.7\*\*\* |
| movement | **1493.8\*\*** | 260.5\*\*\* |
| **unsuccessful** | | |
| eyetracking | 752.1 | **3378.9** |
| no-feedback | 619.5 | **1891.7** |
| movement | 602.6 | **2113.1** |

*listeners spend more time looking at what they think is the referent than at other buttons*

# Interaction Effectiveness

average number of help
requests per interaction

|  | confusion | success |
|---|---|---|
| eyetracking | 1.14 | 91.9 |
| no-feedback | 2.26** | 83.5** |
| movement | 1.77* | 87.5 |

gaze-based feedback makes users
more confident in the interaction

# Interaction Effectiveness

|             | confusion | success |
|-------------|-----------|---------|
| eyetracking | 1.14      | 91.9    |
| no-feedback | 2.26**    | 83.5**  |
| movement    | 1.77*     | 87.5    |

*proportion of correctly resolved referring expressions*

*tracking listener's gaze enhances referential success*

# Interim summary II

- Listeners reliably inspect understood referents in all conditions

- Gaze feedback results in:

  - Lower confusion

  - Positive feedback: speeds interaction

  - Negative feedback: increases success

- (But timing remains an issue!)

# Timing

# Human speaker?

- Is this how human speakers use listener gaze?

- Which eye-movements do they rely on?

- What does feedback really look like?

# Setup

- Walker (12 pairs)
  - Unknown location
  - Eye-tracked by PUPIL P
  - Hears instructions

- Instructor
  - Map
  - Sees walker scene view
  - Gives instructions

# Task

- Walker needs to find table (makro)

- Then walker takes objects and puts them aside (mikro)

- 9 thematic tables with 3-4 target objects each

  - 3 tables in each condition

  - ~40min total

# Conditions

1. GAZE : Walker gaze available to instructor

2. Man-GAZE : Walker gaze perturbed (20% random shift)

3. No-GAZE: Walker gaze NOT available

# Setup



"Could you please pick up the pin box .. eh… that's furthest away from you."

# Measures

Dependent Variables (DV)

1. Instructor behavior:

   a) no. words, feedback (Q1)

2. Low-level listener eye-movements (Q2-i)

3. Listener eye-movements in relation to feedback (Q2-ii)

# Preprocessing DV1

- Transcription

- Forced alignment

- Automatic lemmatization, tagging, shallow syntactic analyses (TreeTagger, Schmid 1995)

- Semi-manual annotation of feedback instances (neg. & pos.)

# Preprocessing DV2

- Standard dispersion-based fixation detection algorithm (Salvucci & Goldberg, 2000)

    - "sequence of gaze points to be a fixation if the maximum distance from their joint center is less than 5% of the scene camera width and the sequence has a minimum duration of 66 msec"

- Sliding window (500ms, step size 250ms) to extract eye movement features, resulting in a dataset of 18,841 time windows

# Preprocessing DV2

| | |
|---|---|
| Fixation | rate, mean, max, variance of durations<br>mean, variance of variance within one fix. |
| Saccades | rate, ratio of (small/large/right/left) sacc.<br>mean, max, variance of amplitudes |
| Combined | ratio saccades / fixations |
| Wordbooks | number of non-zero entries<br>maximum and minimum entries as well as<br>their difference for n-grams with $n \leq 4$ |
| Ratios | all fixation, saccade and combined features<br>in ratio to the value over the whole trial for<br>a particular pair and condition. |

# Preprocessing DV2

- Minimal-redundancy-maximal-relevance criterion (mRMR)

  - Maximizes feature's relevance in terms of mutual information between target variable and features while discarding redundant features (Peng, 2007)

  ➡ **Saccade rate**

# Results - Language

- Performance

  - Success rate  ✗

  - Trial duration  ✗          (Ceiling)

- Language

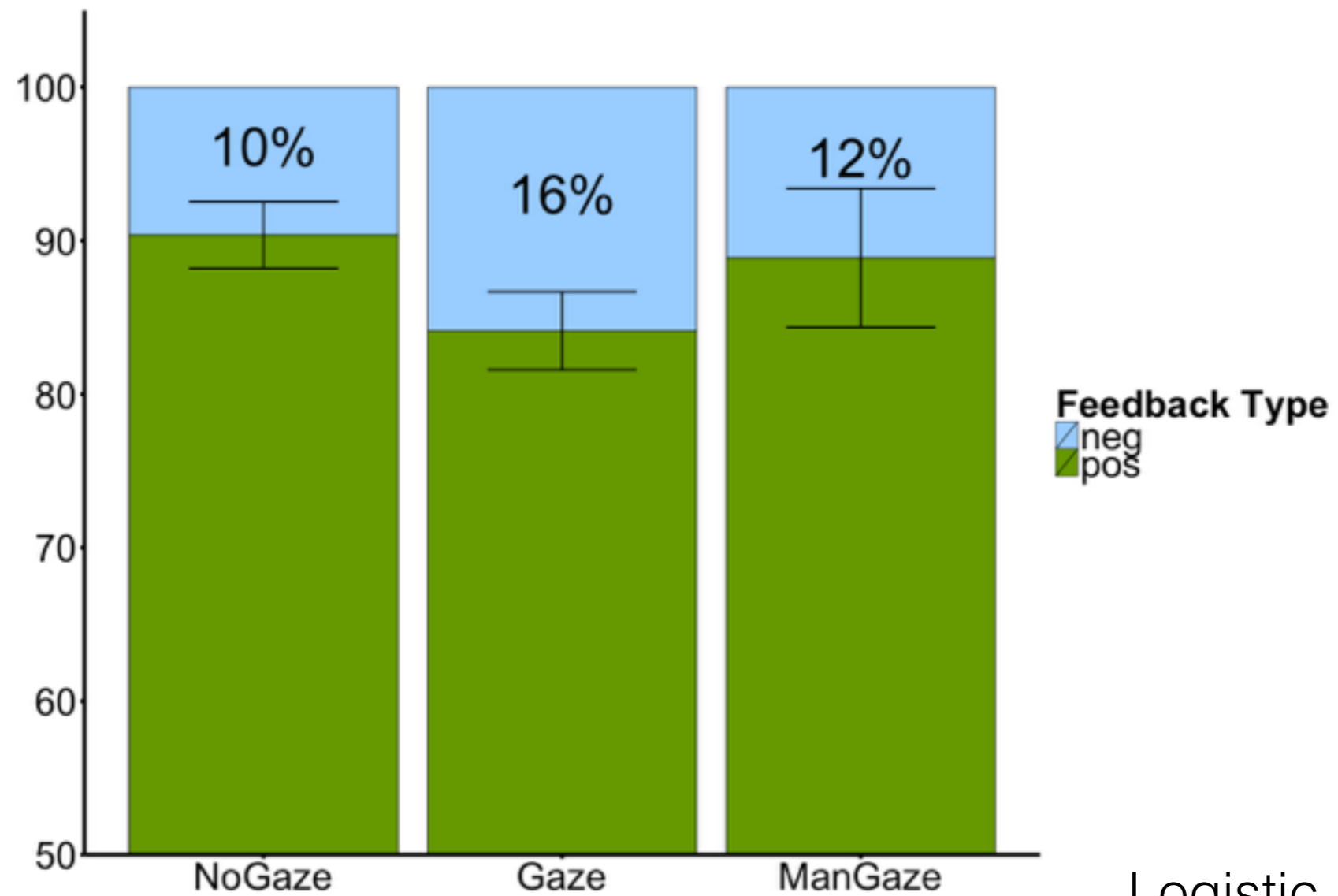  - No. of spoken words  ✗

  - No. of feedback instances  ✓  (Instruction change)
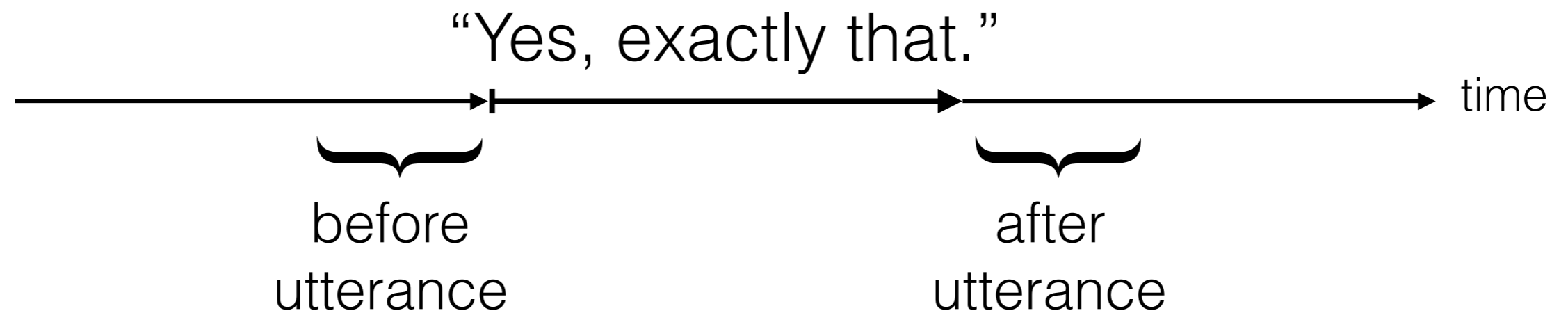
  - Feedback style  ❓

# Results - Feedback



Logistic Regression:
Marg.signifiant

# Low-level eye-movements

"Yes, exactly that."

before utterance

after utterance

time

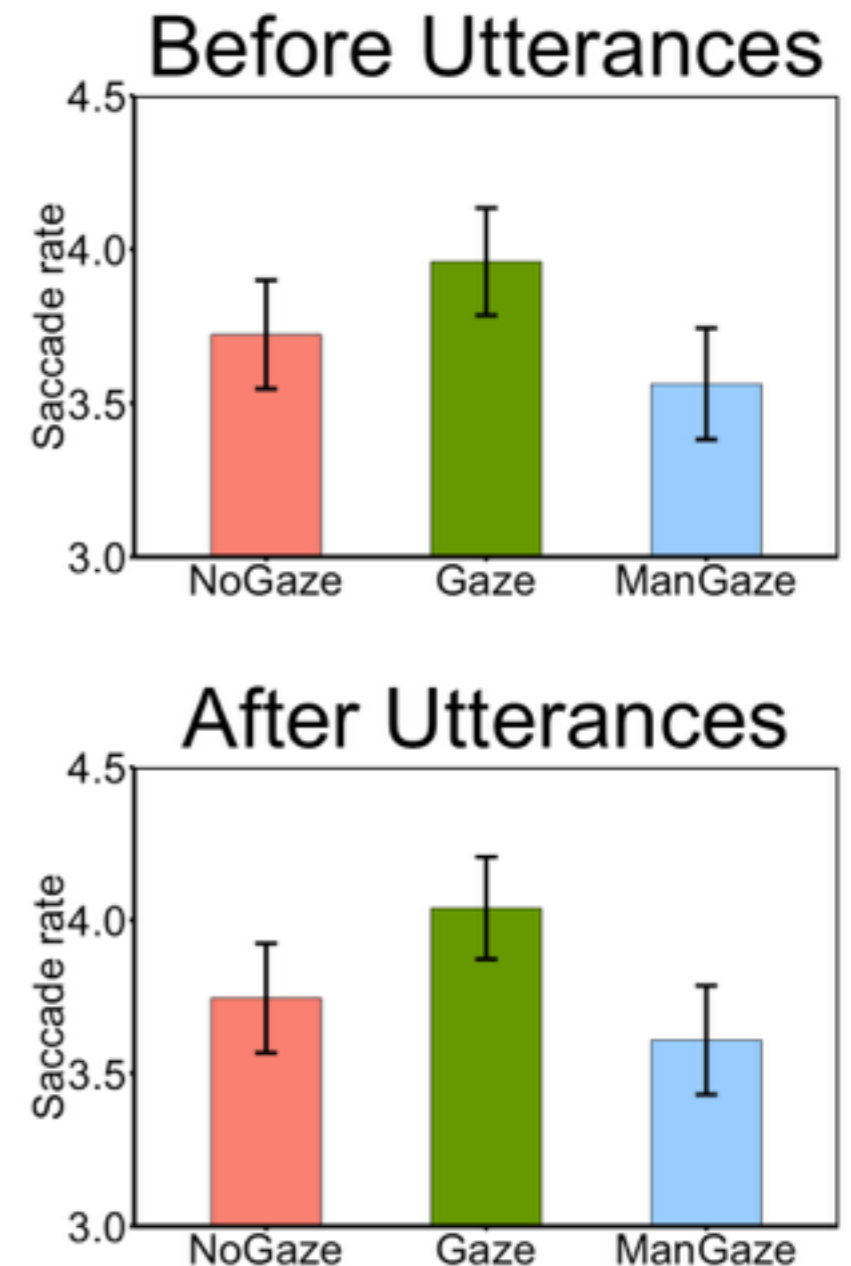Utterance&Presence

# Results - Eye-movements

- Effect of *UtterancePresence* on saccade rate *(task recognition)*

  - No effect of *GazeAvailability*

    **symptom**

- Effect of *FeedbackPresence*

  - Interaction with *GazeAvailability*

- Effect of *GazeAvailability* on saccade rate before & after utterances

    **signal?**



Before Utterances

After Utterances

# Eye-movements & Feedback

- **Manual** annotation of fixations (to target/distractors) up to **5 sec prior to feedback onset**

- No effect of GazeAvailability on patterns found

➡ Feedback timing independent of listener gaze?

# Interim summary III

- Instructions change slightly when listener gaze is available

    - More negative feedback

    - But no measurable effect on performance

- Feedback difficult to categorize

- Eye-movement patterns reflect speech process **symptom & signal** AND change with GazeAvailability

- Gaze-Feedback pattern constant across conditions

# Conclusion

- Listeners follow speaker gaze (and arrows) and form predictions about upcoming referents

  - Difference in strategic use of these cues

- Listeners follow speech & these gaze cues can be exploited by the speaker

  - Efficient use by NLG system

  - Little benefit for real speaker

    - Ceiling, Unnatural situation

# References

- Hanna, J. & Brennan, S. (2007). JML

- Staudte, M. & Crocker, M. (2011). Cognition.

- Staudte, M., Crocker, M., Heloir, A., & Kipp, M. (2014). Cognition

- Garoufi, K., Staudte, M., Koller, A., & Crocker, M. (2015). Cognitive Science, in press.

- Koleva, N., Hoppe, S., Staudte, M., & Bulling, A. (2015). Proc. of the Annual Meeting of the Cognitive Science Society, Los Angeles.