



How Do We Learn about some New: Interpreting Generic Statements Using Bayesian Inference

Just in case you don't
know it already:

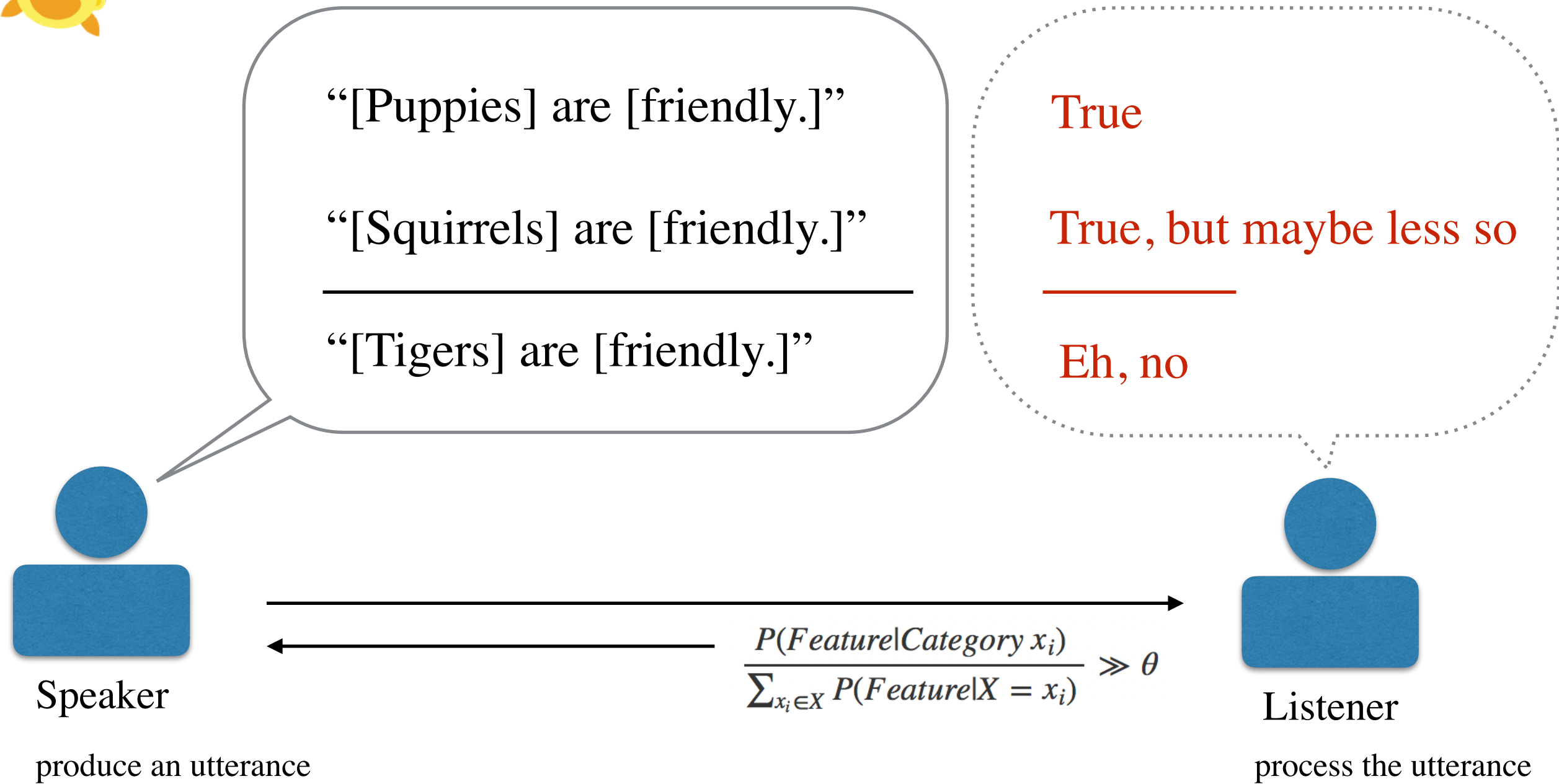
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Presenter: Flora Zhang

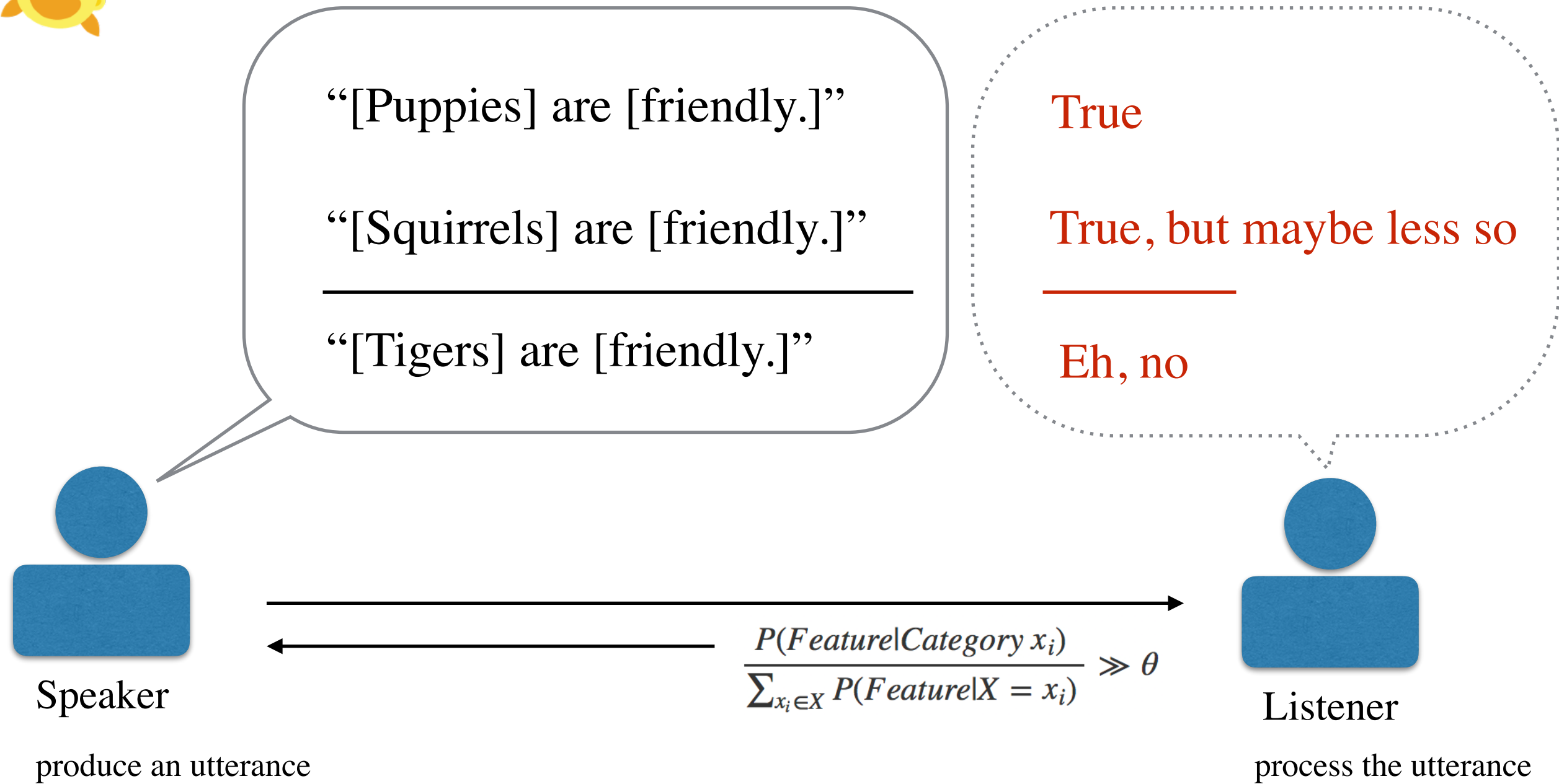
Principle Investigator: Dr. Daniel Yurovsky

MACS Github Repo: <https://github.com/xiuyuanzhang/MACS30200proj>

Research Project Repo: <https://github.com/xiuyuanzhang/generic-statement>



Informally defined, generic statements are blanket statements about members of a category and their diagnostic features.




Informally defined, generic statements are blanket statements about members of a category and their diagnostic features.

“[Puppies] are [friendly.]”

“[Squirrels] are [friendly.]”

“[Tigers] are [friendly.]”

“feps” is a made-up word,
a novel category



THE UNIVERSITY OF
CHICAGO

DIVISION OF THE SOCIAL SCIENCES


The Akarians tell you that feps are like goats.

One of the Akarians says that: "feps are friendly."

What percent of feps do you think are friendly?

You can drag the slider bar below to show your results.

0% 100%





A Speaker-Listener Interaction Model

$$\frac{P(\text{Feature}|\text{Category } x_i)}{\sum_{x_i \in X} P(\text{Feature}|X = x_i)} \gg \theta$$

$$\frac{P(\text{Feature}|\text{Novel Category})}{\sum_{x_i \in X} P(\text{Feature}|X = x_i)} \gg \theta$$

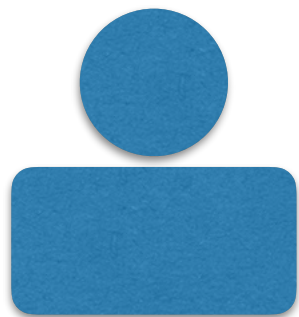
$$P(\text{Feps are friendly}|\text{speaker uttered "Feps are friendly"})$$

$$= \frac{P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})P(\text{Feps are friendly})}{P(\text{speaker uttered "Feps are friendly"})}$$

$$= \mathbb{E}(P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})P(\text{Feps are friendly}))$$

$$= \mathbb{E}(P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})) \cdot 0.5$$

$$= \mathbb{E}\left(\int_{P(\text{Friendly}|\text{Feps})} \frac{e^{\alpha P(\text{Friendly}|\text{Feps})}}{e^{\alpha P(\text{Friendly}|\text{Feps})} + e^{\alpha P(\text{Friendly})}} \delta P(\text{Friendly}|\text{Feps})) \cdot 0.5\right)$$



Listener

process the utterance



A Speaker-Listener Interaction Model

$$\frac{P(\text{Feature}|\text{Category } x_i)}{\sum_{x_i \in X} P(\text{Feature}|X = x_i)} \gg \theta$$

$$\frac{P(\text{Feature}|\text{Novel Category})}{\sum_{x_i \in X} P(\text{Feature}|X = x_i)} \gg \theta$$

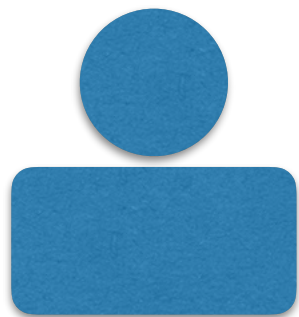
$$P(\text{Feps are friendly}|\text{speaker uttered "Feps are friendly"})$$

$$= \frac{P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})P(\text{Feps are friendly})}{P(\text{speaker uttered "Feps are friendly"})}$$

$$= \mathbb{E}(P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})P(\text{Feps are friendly}))$$

$$= \mathbb{E}(P(\text{speaker uttered "Feps are friendly"}|\text{Feps are friendly})) \cdot 0.5$$

$$= \mathbb{E}\left(\int_{P(\text{Friendly}|\text{Feps})} \frac{e^{\alpha P(\text{Friendly}|\text{Feps})}}{e^{\alpha P(\text{Friendly}|\text{Feps})} + e^{\alpha P(\text{Friendly})}} \delta P(\text{Friendly}|\text{Feps})) \cdot 0.5\right)$$



Listener

process the utterance



Method

Data:

- Involves Human Subjects (SBS IRB Approval No.: IRB16-1118) - Online Survey using Amazon MTurk
- Pew Research Center's American Trends Panel Wave 26 (n = 5155)

Analysis and Modeling in R & Stan:

- Random Forest - finding distinct demographic trends on social and political issues (for next step surveys)
- Linear Regression - evaluating estimate probability for novel category
- Logistic Regression - evaluating binary True/False var
- Stan(probabilistic programming language) - build Bayesian statistical model for our hypothesis

Previous Related Work

- Ward, Andrew, L. Ross, E. Reed, E. Turiel, and T. Brown. "Naive realism in everyday life: Implications for social conflict and misunderstanding." *Values and knowledge* (1997): 103-135.
- Rhodes, Marjorie, Sarah-Jane Leslie, and Christina M. Tworek. "Cultural transmission of social essentialism." *Proceedings of the National Academy of Sciences* 109, no. 34 (2012): 13526-13531.
- Tessler, Michael Henry, and Noah D. Goodman. "A pragmatic theory of generic language." *arXiv preprint arXiv:1608.02926*(2016).



Method

Data:

- Involves Human Subjects (SBS IRB Approval No.: IRB16-1118) - Online Survey using Amazon MTurk
- Pew Research Center's American Trends Panel Wave 26 (n = 5155)

Analysis and Modeling in R & Stan:

- Random Forest - finding distinct demographic trends on social and political issues (for next step surveys)
- Linear Regression - evaluating estimate probability for novel category
- Logistic Regression - evaluating binary True/False var
- Stan(probabilistic programming language) - build Bayesian statistical model for our hypothesis

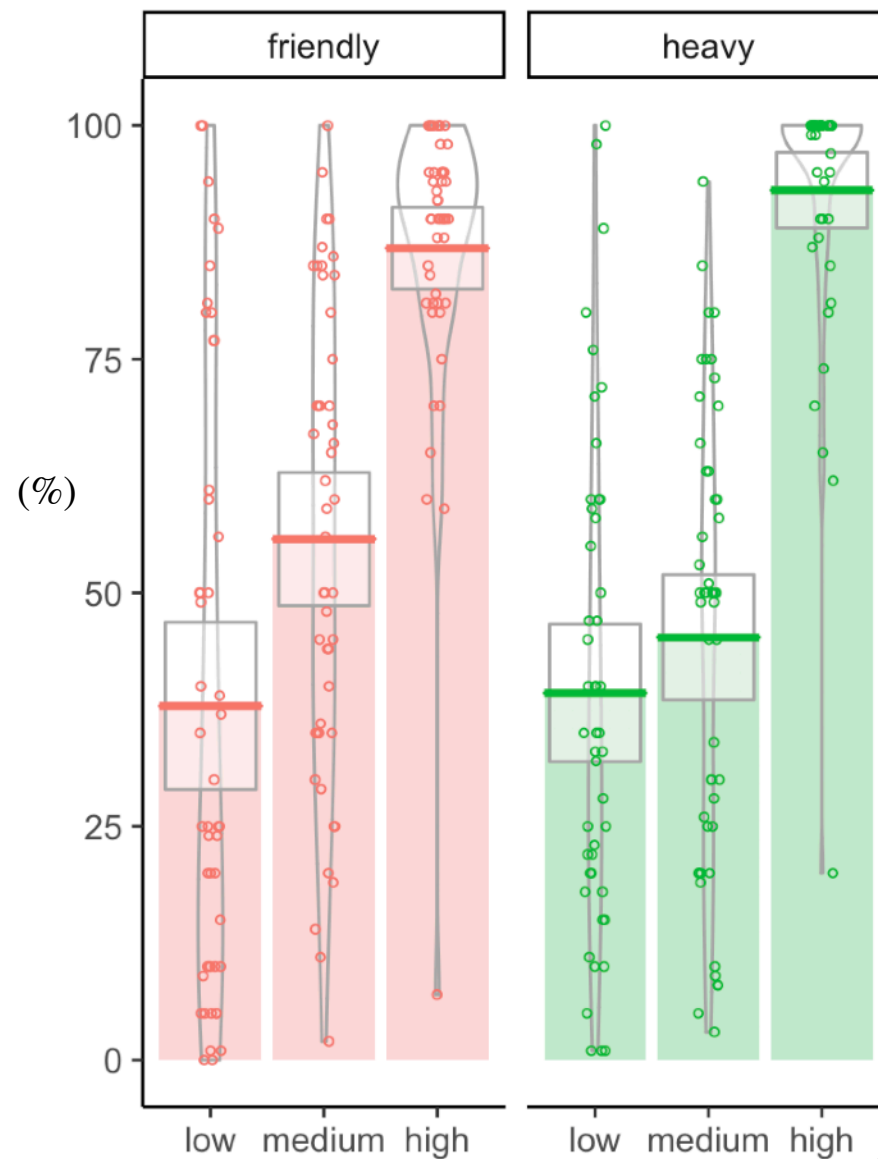
Previous Related Work

- Ward, Andrew, L. Ross, E. Reed, E. Turiel, and T. Brown. "Naive realism in everyday life: Implications for social conflict and misunderstanding." *Values and knowledge* (1997): 103-135.
- Rhodes, Marjorie, Sarah-Jane Leslie, and Christina M. Tworek. "Cultural transmission of social essentialism." *Proceedings of the National Academy of Sciences* 109, no. 34 (2012): 13526-13531.
- Tessler, Michael Henry, and Noah D. Goodman. "A pragmatic theory of generic language." *arXiv preprint arXiv:1608.02926*(2016).



Preliminary Results

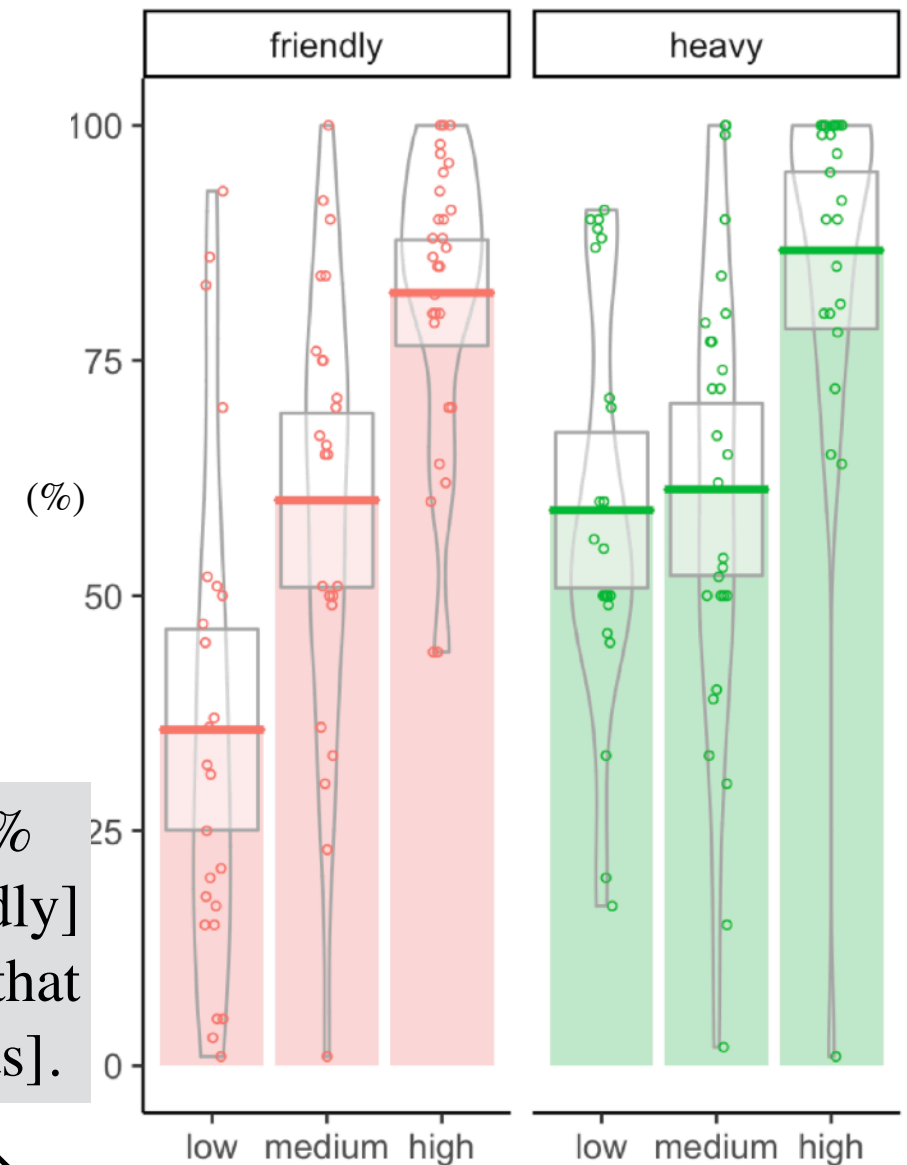
Baseline Generic Estimation by Features (n = 145)



Participants say __% of
[goats] are [friendly].

Participants say __%
of [feps] are [friendly]
when we tell them that
[feps] are like [goats].

Novel Generic Estimation by Features (n = 78)



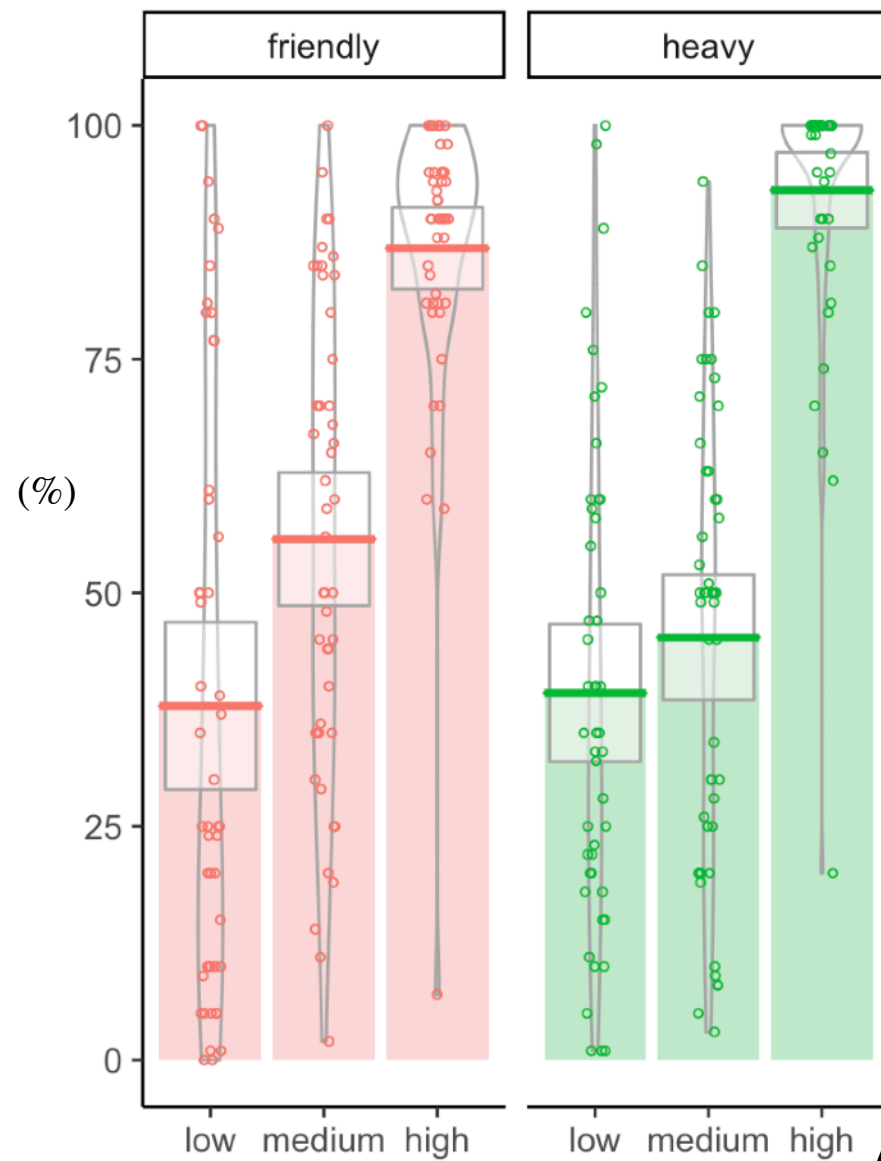
Experiment Conditions for Referenced Familiar Categories by Levels of Probability Intensity

Survey participants: n = 300 (150 each)



Preliminary Results

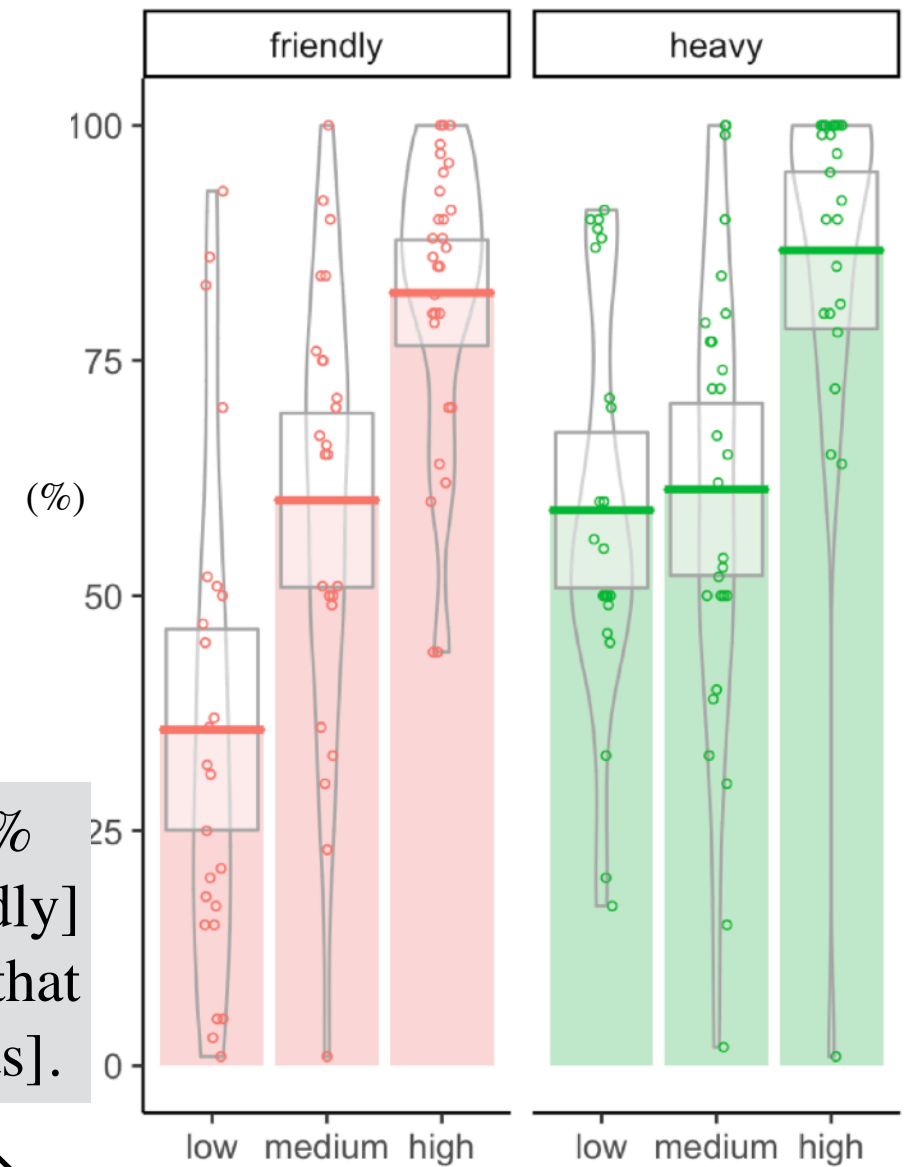
Baseline Generic Estimation by Features (n = 145)



Participants say __% of
[goats] are [friendly].

Participants say __%
of [feps] are [friendly]
when we tell them that
[feps] are like [goats].

Novel Generic Estimation by Features (n = 78)



Experiment Conditions for Referenced Familiar Categories by Levels of Probability Intensity

Survey participants: n = 300 (150 each)