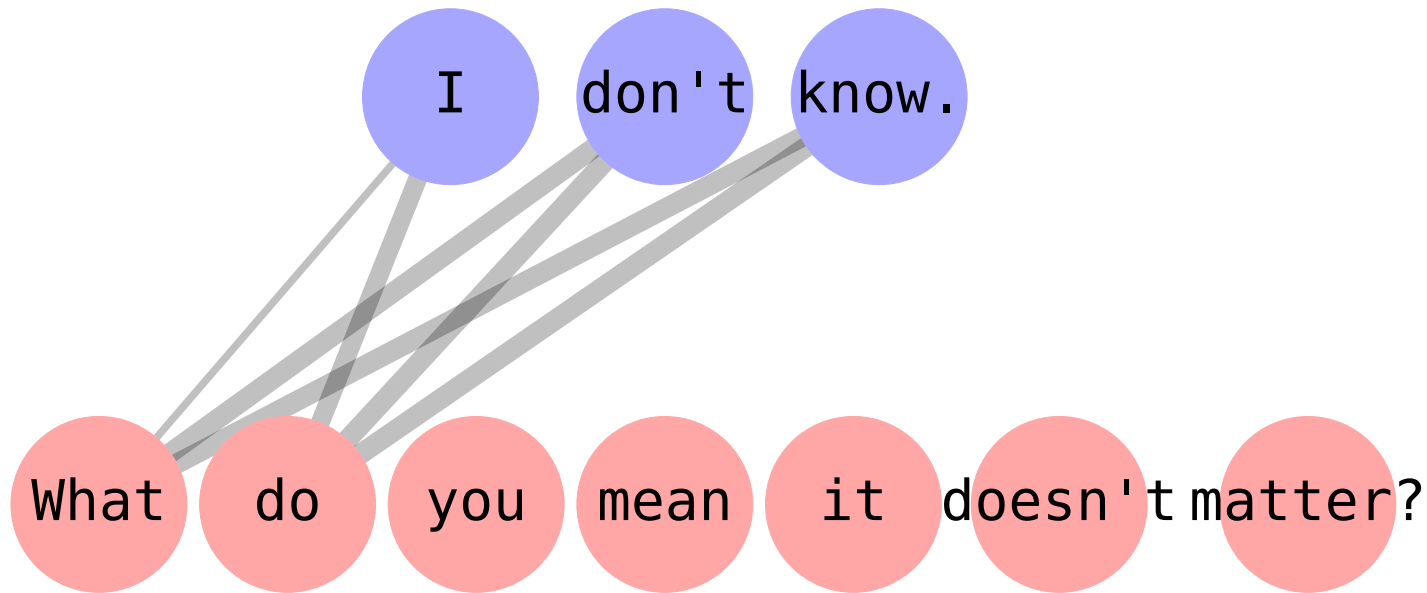
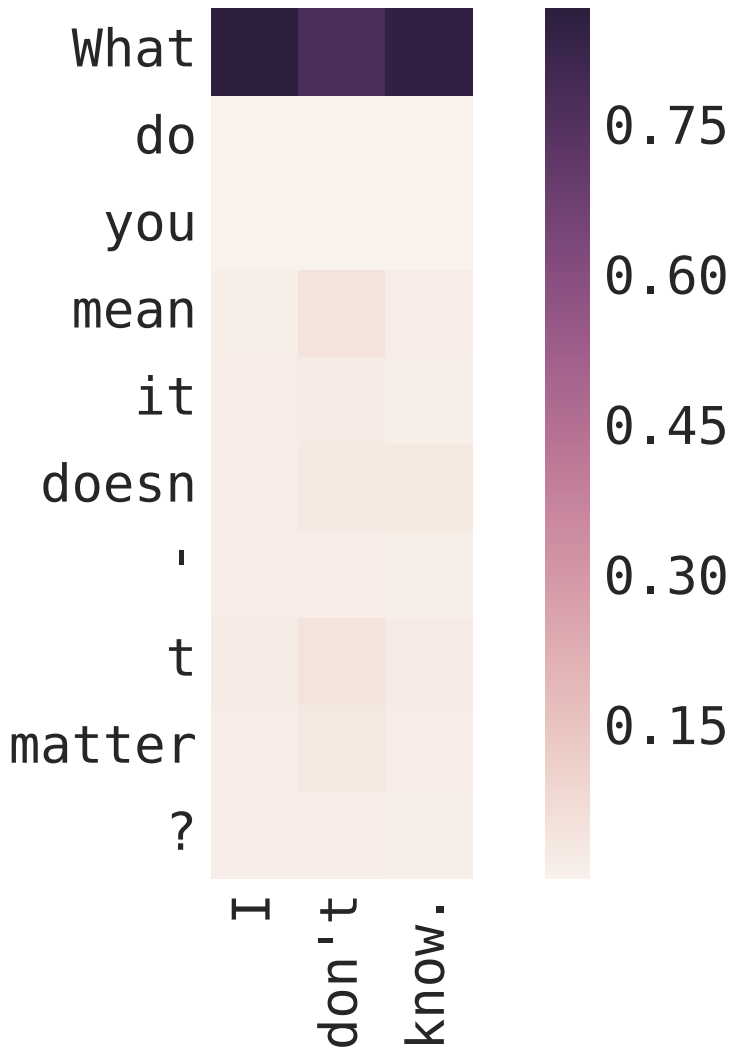


• Black-box: seq2seq + attention, trained on dialogue corpus

Interpretability for Flaw Detection in Dialogue Systems







*What do you mean
It doesn't matter?*

Input



Neural Net



I don't know

Prediction

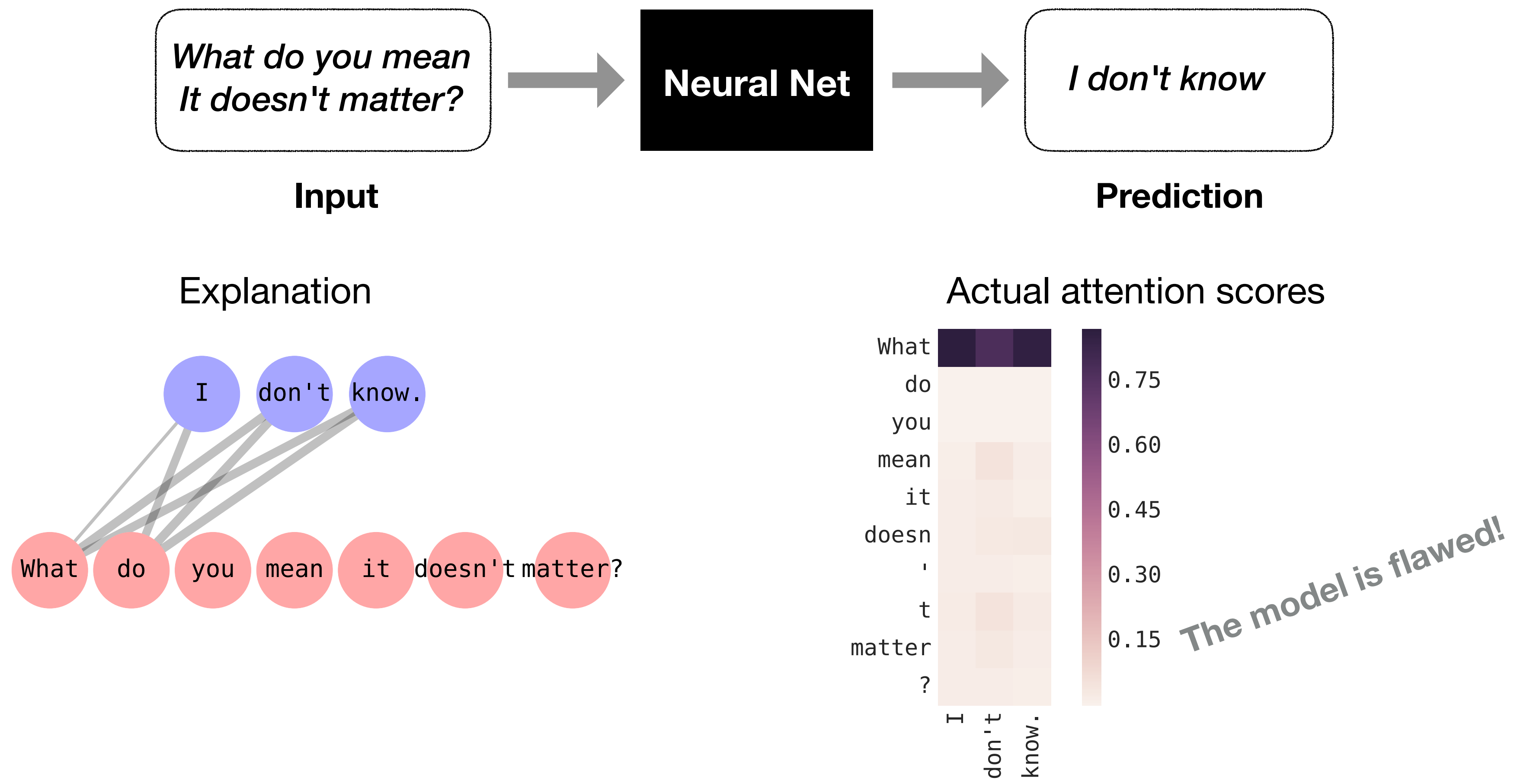
The model is flawed!

Explanation

Actual attention scores

Interpretability for Flaw Detection in Dialogue Systems

- Black-box: seq2seq + attention, trained on dialogue corpus



The only option for
already-trained models

Works for
Black-Box models



Post-hoc Interpretability