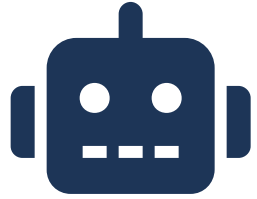


Harnessing explanations to bridge AI & humans



Vivian Lai, Samuel Carton, Chenhao Tan
@vivwylai | @SamHCarton | @ChenhaoTan
University of Colorado Boulder
CHI2020 Fair & Responsible AI Workshop

Ubiquitous Machine Learning



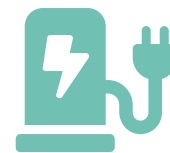
Medical diagnosis



Credit score prediction

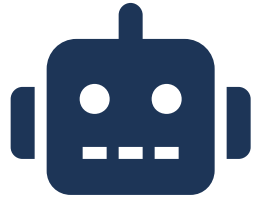


Recidivism prediction



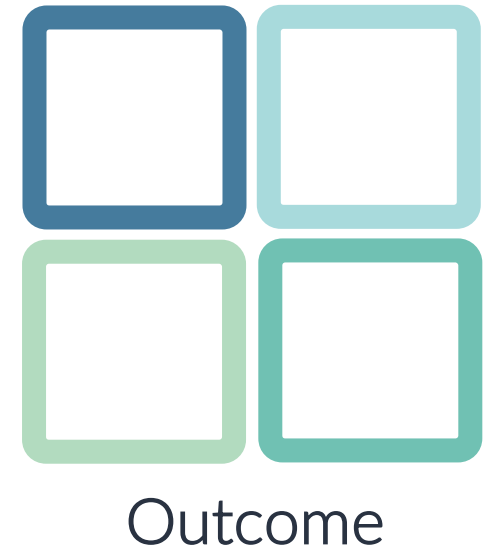
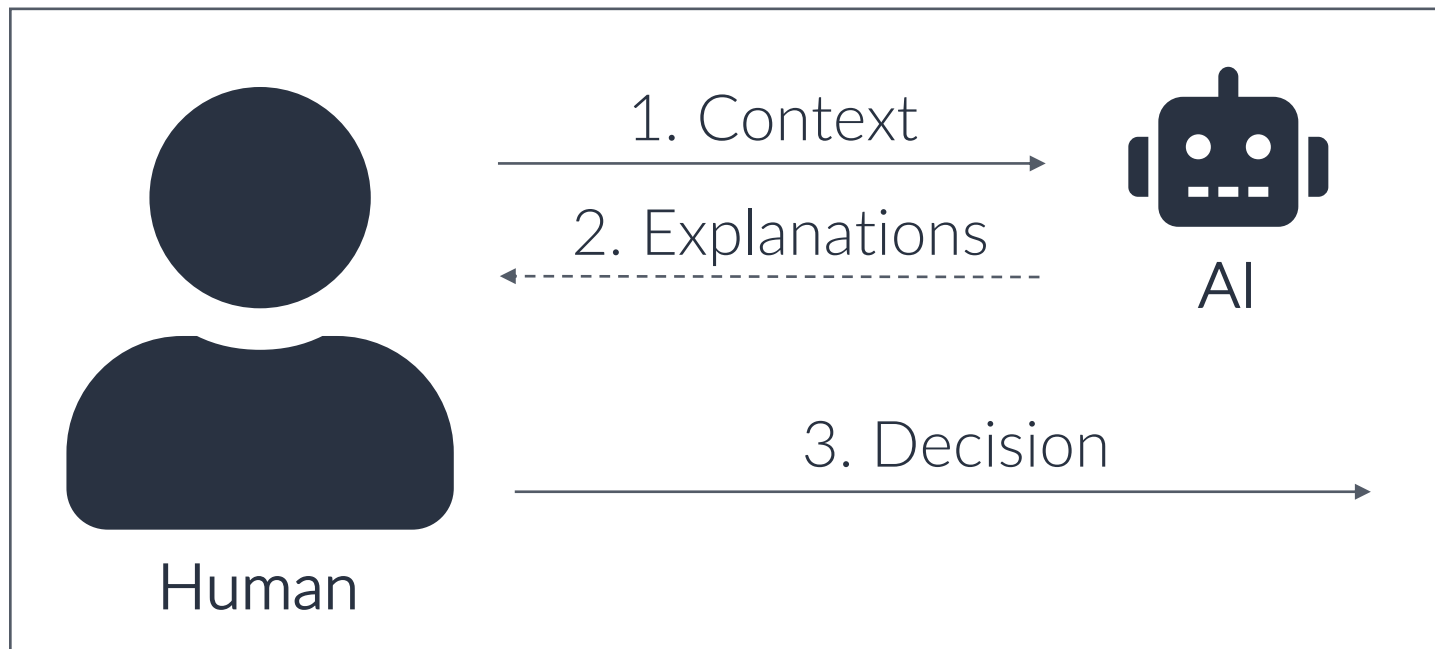
Autonomous driving

Full automation is not desired

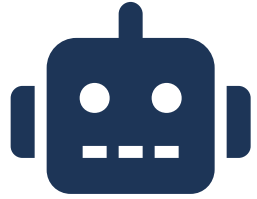


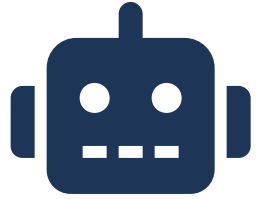
Machine-in-the-loop decision making

Challenging tasks



Fail to improve performance





Fail to improve performance



You made **Chicago** a wonderful stay! The room was gorgeous! I came with **very** little **on** hand and **my** deluxe room supplied me with everything **that** I needed, I didn't even have to **ask**! Thank you so much, **I will be back**! **Very** tidy room as well!

You made **Chicago** a wonderful stay! The room was gorgeous! I came with **very** little **on** hand and **my** deluxe room supplied me with everything **that** I needed, I didn't even have to **ask**! Thank you so much, **I will be back**! **Very** tidy room as well!

AI can discover inconspicuous and counterintuitive patterns



Misalignment between explanations and human mental model

Augment human
mental model

Towards human-
centered
explanations



Augment human mental model



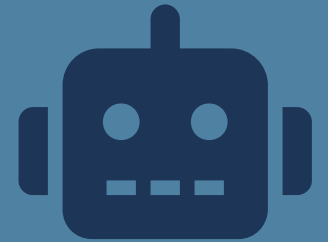
Model-driven tutorials

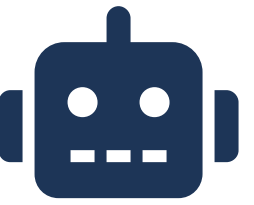


Interactive explanations

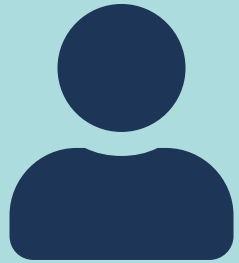


Evaluating generalization





Towards human-centered explanations



Understanding human explanations



Experimenting with alternative explanation types



Explanations as model criticism

Vivian Lai, Samuel Carton, Chenhao Tan
@vivwylai | vivwylai@gmail.com
@SamHCarton | @ChenhaoTan
University of Colorado Boulder



workshop: <https://tinyurl.com/harness-explanations>
full paper: <https://tinyurl.com/model-driven-tutorials>