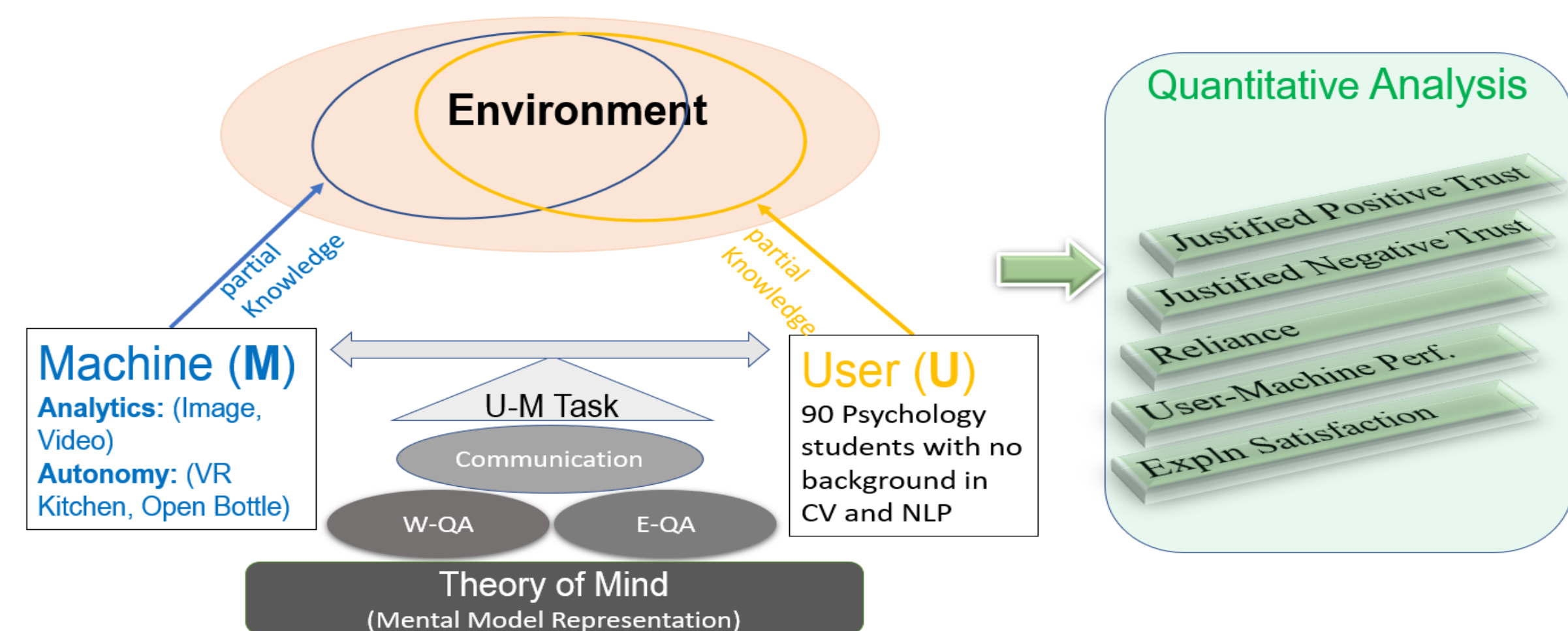
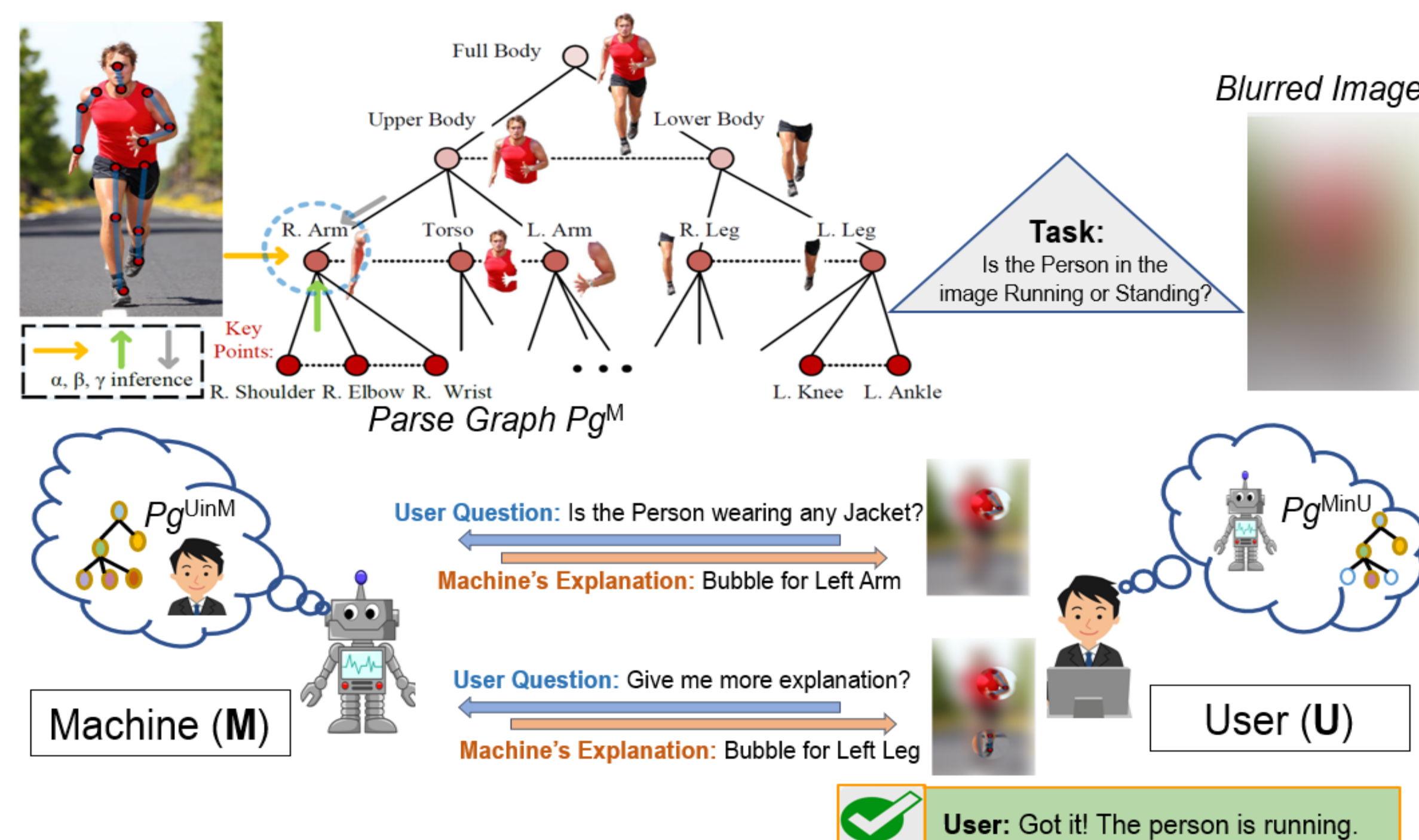


X-ToM: EXPLANATION WITH THEORY-OF-MIND



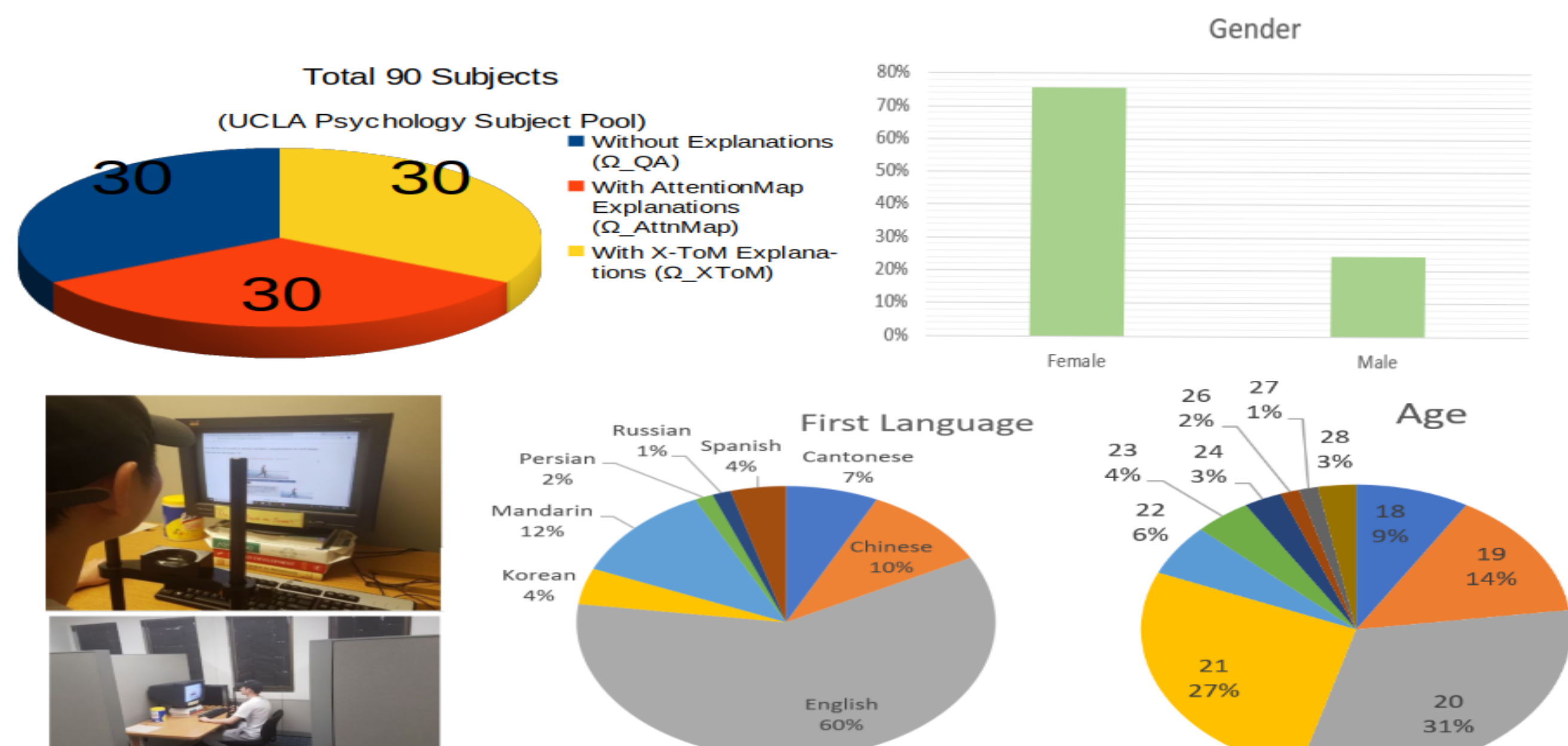
X-ToM for optimizing the dialog with a user towards estimating and increasing human trust.

XAI AS COLLABORATIVE TASK SOLVING

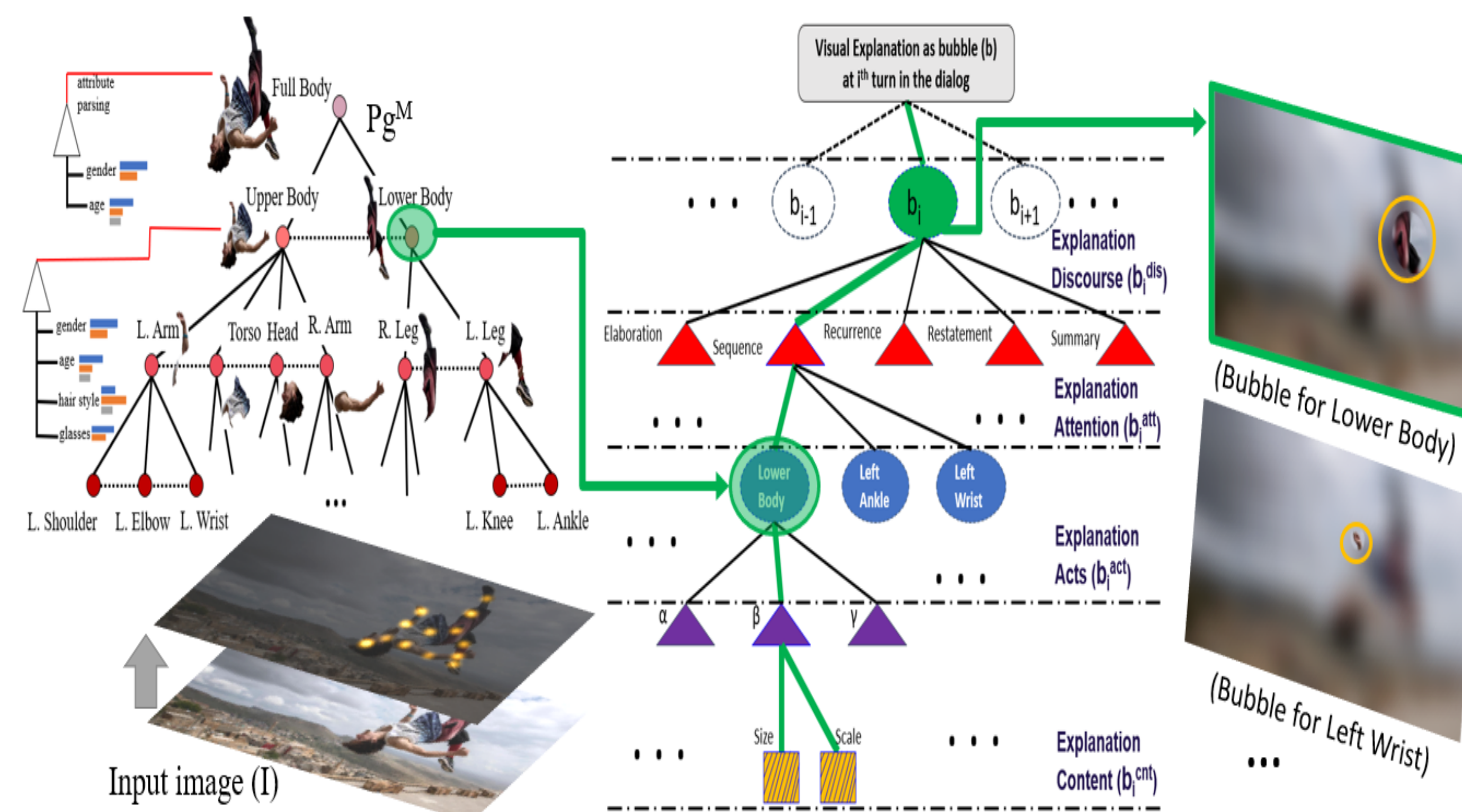


Through a dialog, we estimate Trust and Reliance in terms of pg^M , pg^{UinM} and pg^{MinU} .

SUBJECT POOL



EXPLANATION GENERATION IN X-ToM



Both the Machine and the User solve the image recognition tasks. The Machine interprets the image I as pg^M . The Human receives visual explanations – bubbles – optimized by the X-ToM Explainer.

TRUST ESTIMATION

It is possible for humans to feel positive trust with respect to certain tasks, while feeling negative trust (i.e. mistrust) on some other tasks. **Justified Positive Trust (JPT):**

$$JPT = \frac{1}{N} \sum_i \sum_{z=\alpha, \beta, \gamma} \Delta JPT(i, z),$$

$$\Delta JPT(i, z) = \frac{\|pg_{i,z,+}^{MinU} \cap pg_{i,+}^M\|}{\|pg_{i,+}^M\|}$$

Justified Negative Trust (JNT):

$$JNT = \frac{1}{N} \sum_i \sum_{z=\alpha, \beta, \gamma} \Delta JNT(i, z),$$

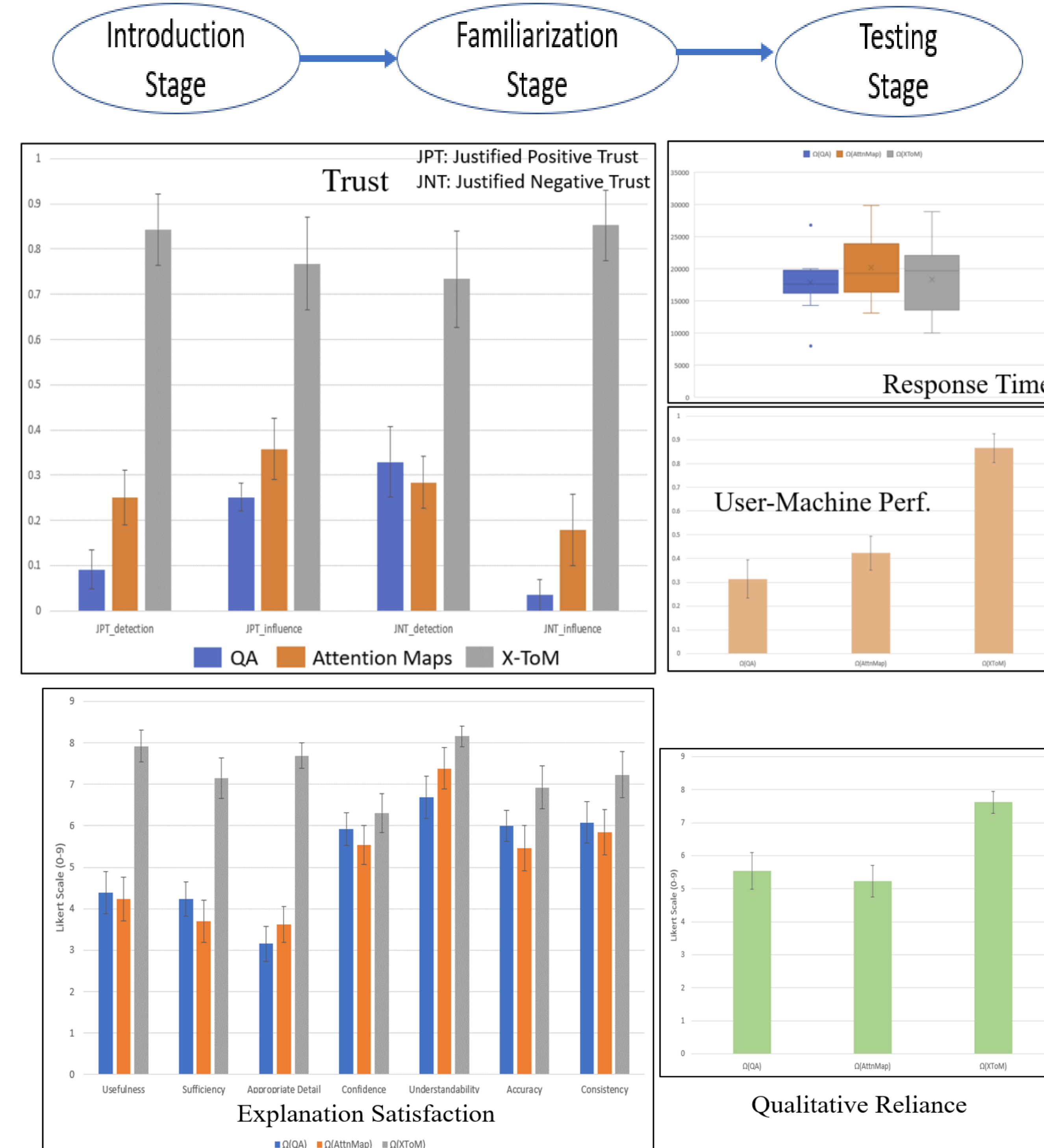
$$\Delta JNT(i, z) = \frac{\|pg_{i,z,-}^{MinU} \cap pg_{i,-}^M\|}{\|pg_{i,-}^M\|}$$

Reliance (Rc): Reliance captures the extent to which a human can accurately predict the performer's inference results without over- or under-estimation.

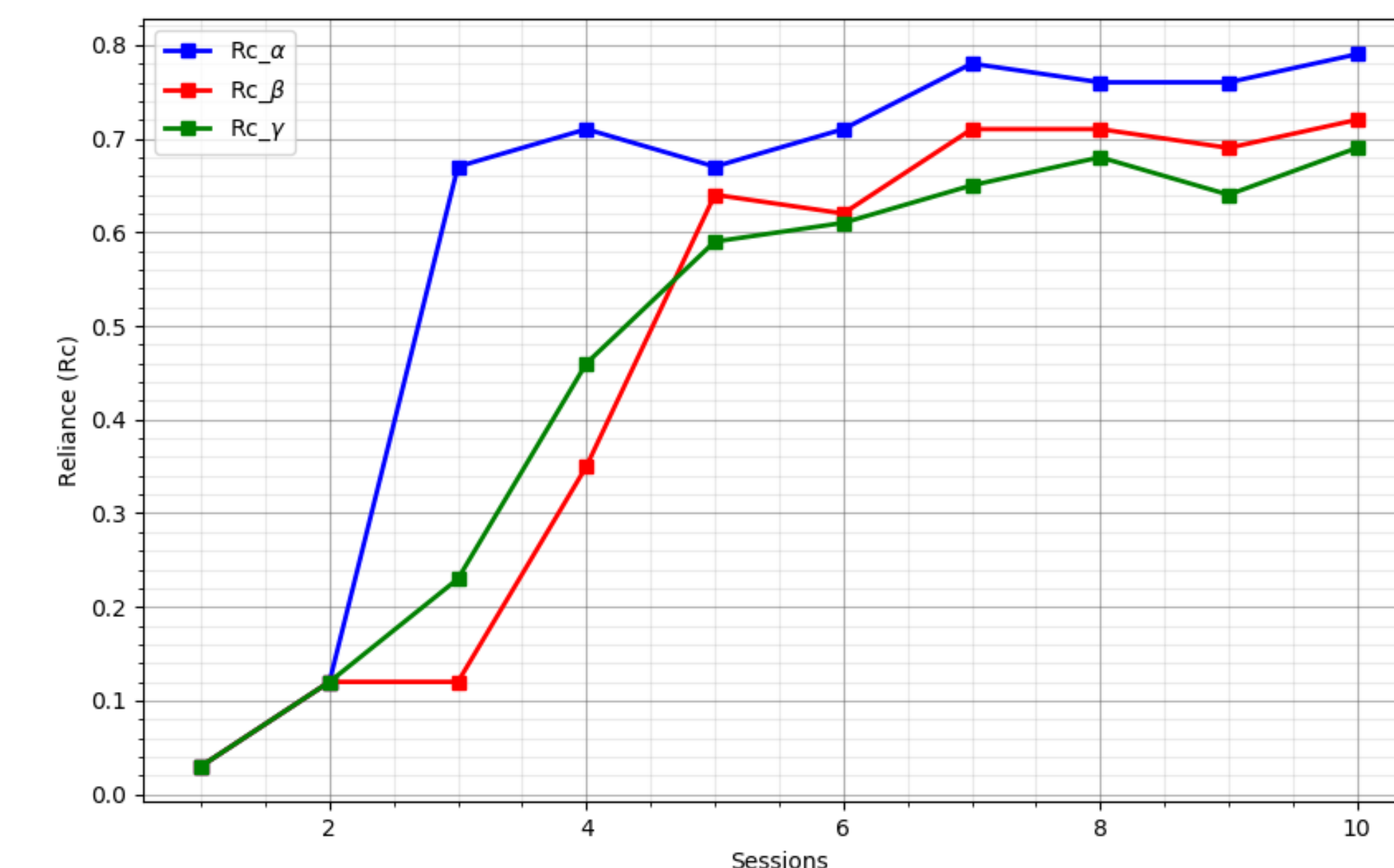
$$Rc = \frac{1}{N} \sum_i \sum_{z=\alpha, \beta, \gamma} \Delta Rc(i, z),$$

$$\Delta Rc(i, z) = \frac{\|pg_{i,z}^{MinU} \cap pg_{i,z}^M\|}{\|pg_{i,z}^M\|}$$

RESULTS OF OUR HUMAN STUDY



X-ToM significantly outperformed ($p < 0.01$) baselines (QA, Attention Maps) in terms of Appropriate Trust, Reliance, User-Machine Performance and Satisfaction.



Gain in Reliance over sessions w.r.t α , β and γ processes. As we can see, with more sessions, we can further improve human reliance.