# "Good Enough" Agents: Investigating Designed Imperfections in Human-AI Teams Across Parallel Domains

**Sebastian S. Rodriguez**
**Doctoral Final Defense**

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN

# Evolving automation



Simple

Complex

Automation Level

Autonomy

Level 1

Level 10

Information Acquisition    **Monitoring**
Information Analysis       **Generating**
Decision Making            **Selecting**
Action Implementation      **Executing**

# Human-agent interaction

Human Teams

HABA-MABA Teams
Function Allocation

Human-Agent Teams
Interdependent Activity

Inter-predictable
Commonly Grounded
Directable

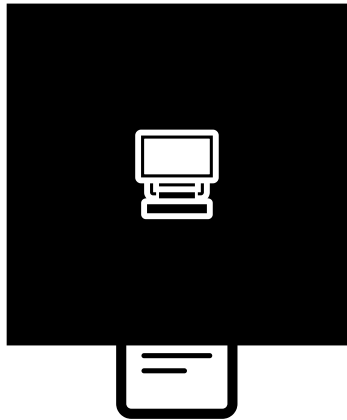# Inappropriate use: we're bad at calibrating

- Trust is a facet researched often w.r.t automation

- Truth is, humans are **bad** at calibrating trust[a]
  - We expect too much
  - We expect too little
  - We can't make up our minds

- Hitting the sweet spot is challenging: calibration like human-human trust[b]

a V. L. Pop, A. Shrewsbury, and F. T. Durso. 2015. Individual differences in the calibration of trust in automation. Human Factors 57, 4, Jun 2015.
b J. D. Lee and K. A. See. 2004. Trust in Automation: Designing for Appropriate Reliance. Human Factors: The Journal of the Human Factors and Ergonomics Society 46, 1, Jan 2004.

# Under-trust has been addressed



Explanations



Control settings
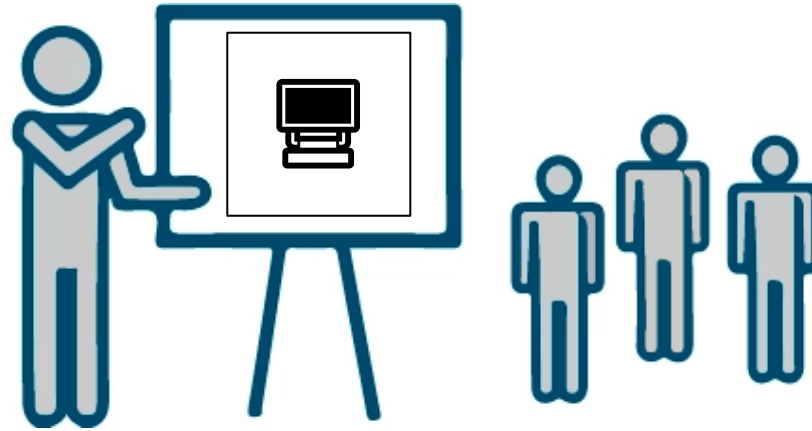


Information displays



Anthropomorphism



UX design

Trust is shown to positively influence technology acceptance[a]
But could these features[b] generate excessive trust?

a K. Wu, Y. Zhao, Q. Zhu, X. Tan, H. Zheng. A meta-analysis of the impact of trust on technology acceptance mode: Investigation of moderating influence of subject and context type. International Journal of Information Management 31, 6, Dec 2011.
b H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, J. Vaughan. Interpreting interpretability: Understanding data scientists' use of interpretability tools for machine learning. CHI 2020 Proceedings. 2020.

# Over-trust is being addressed

**Viable approaches: training[a] and repeated exposure[b] to manual control**
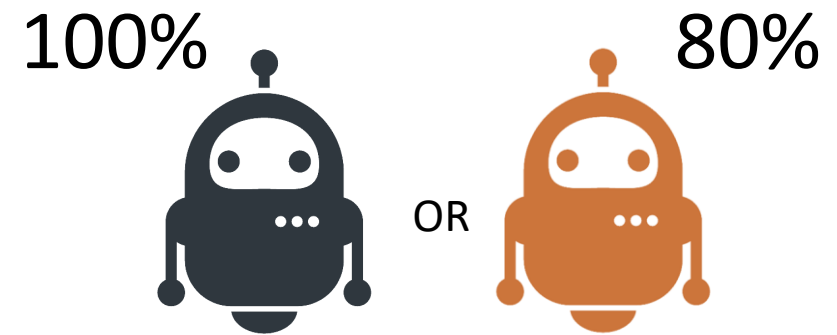Requires time and effort to learn and become skilled at calibrated trust
Any other approaches?

**a** J. D. Lee and K. A. See. 2004. Trust in Automation: Designing for Appropriate Reliance. Human Factors: The Journal of the Human Factors and Ergonomics Society 46, 1, Jan 2004.
**b** M. Itoh. Necessity of supporting situation understanding to prevent over-trust in automation. International Electronic Journal of Nuclear Safety and Simulation 43, 46, 2010.

# "Good Enough"-ness

- We expect AI to be optimal[a]

- What if it doesn't need to be?

- Reliability is often a feature not manipulated as part of agent design
  - Research states AI reliability should be >70% for any benefit[b]

- At which reliability is the most benefit found?
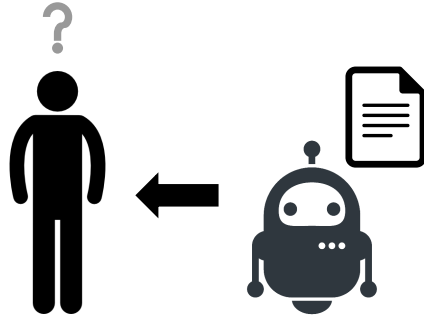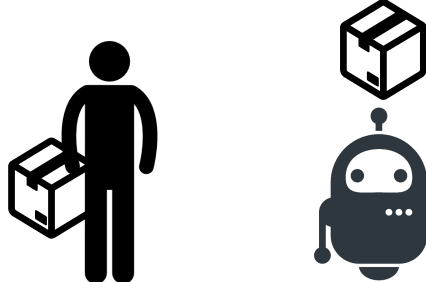
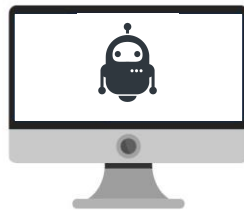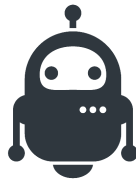- Design goal or minimum returns?

100%  OR  80%

**"Good Enough" Agent**
An agent engrained with small error, to simulate less than perfect performance

a M. Endsley. From Here to Autonomy: Lessons Learned From Human-Automation Research. Human Factors, 59, 1. Feb 2017.
b B. Schelble, C. Flathmann, N. McNeese. Towards Meaningfully Integrating Human-Autonomy Teaming in Applied Settings. International Conference on Human-Agent Interaction, Nov 2020.

# Domains of operation

**a** C. Esterwood, K. Essenmacher, H. Yang, F. Zeng, L. Robert. A Meta-Analysis of Human Personality and Robot Acceptance in Human-Robot Interaction. Conference on Human Factors in Computing Systems. 2021.
**b** J. Li. The benefit of being physically present: A survey of experimental works comparing copresent robots, telepresent robots and virtual agents. International Journal of Human-Computer Studies. 2014.
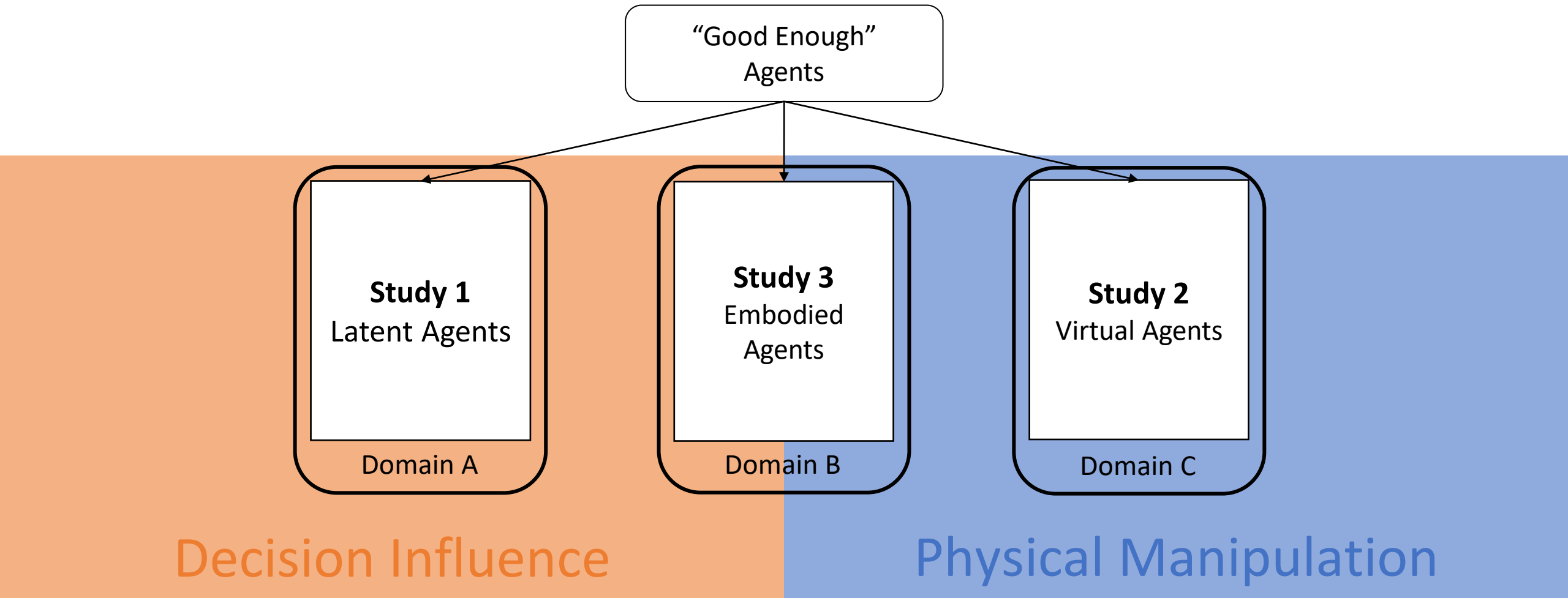
# Research Questions

How does lowering agent reliability (i.e., "Good Enough" agents) affect task performance and human trust across different domains?
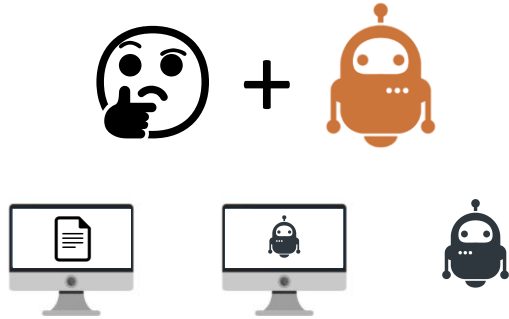
What is the relationship between human individual differences, agent reliability, and technology use?

How do different task requirements interact with the human-AI interaction dynamic?
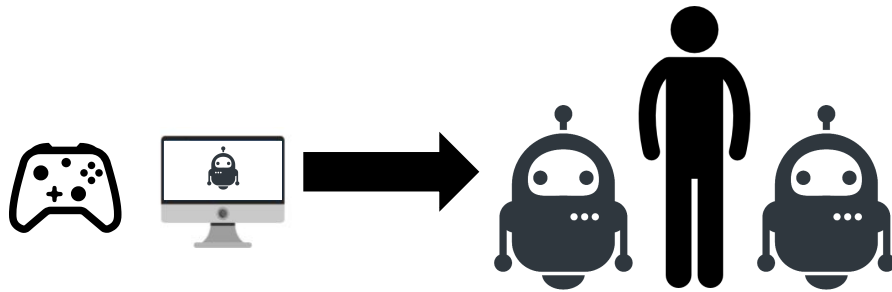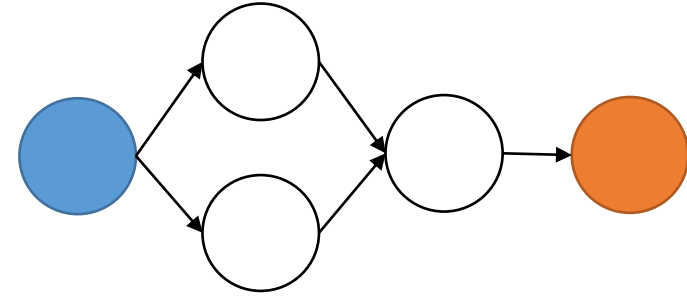
# Dissertation Contributions
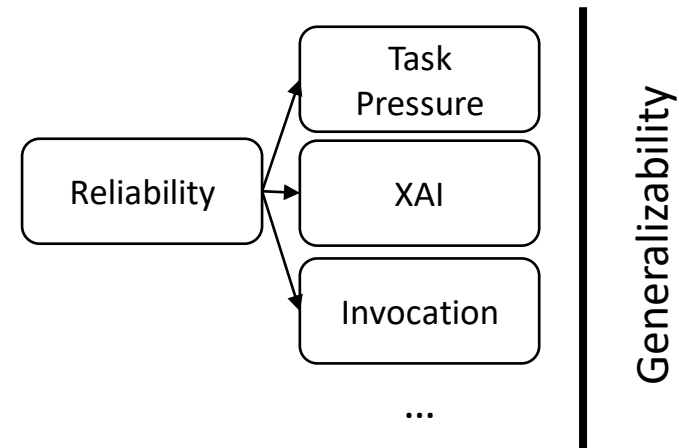


Evidence showing "Good Enough" agents lead to positive outcomes in different contexts



Validation of 3 game-based tasks for human-AI interaction research



Holistic mediation models for human-AI interaction



Reliability → Task Pressure

Reliability → XAI

Reliability → Invocation

...

Generalizability

An ontology of task features that impede generalizability in reliability research
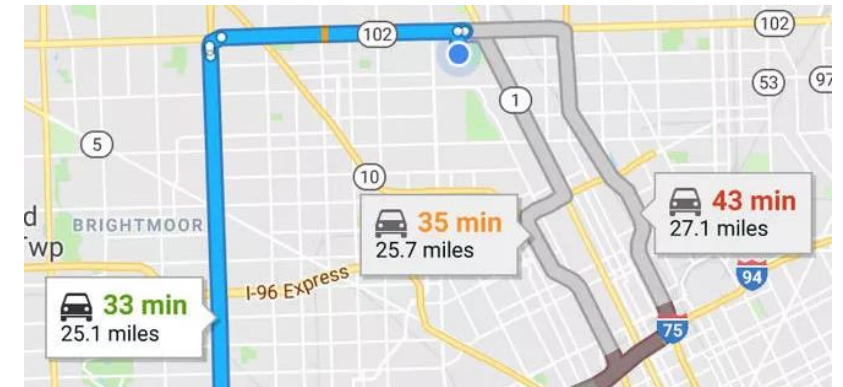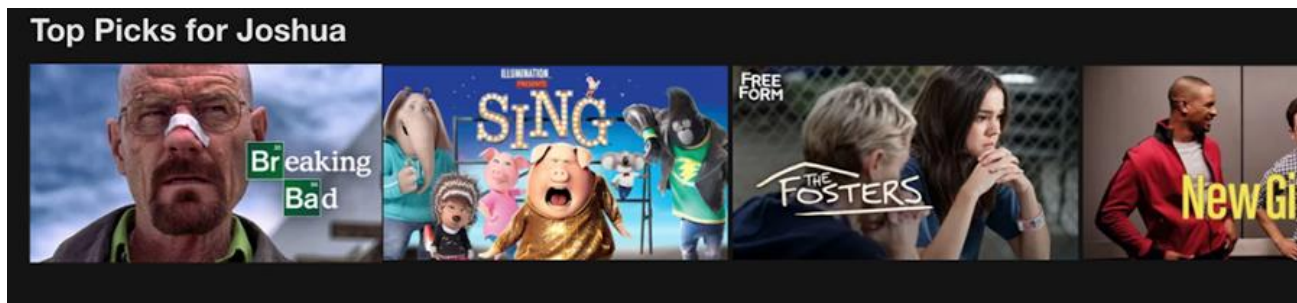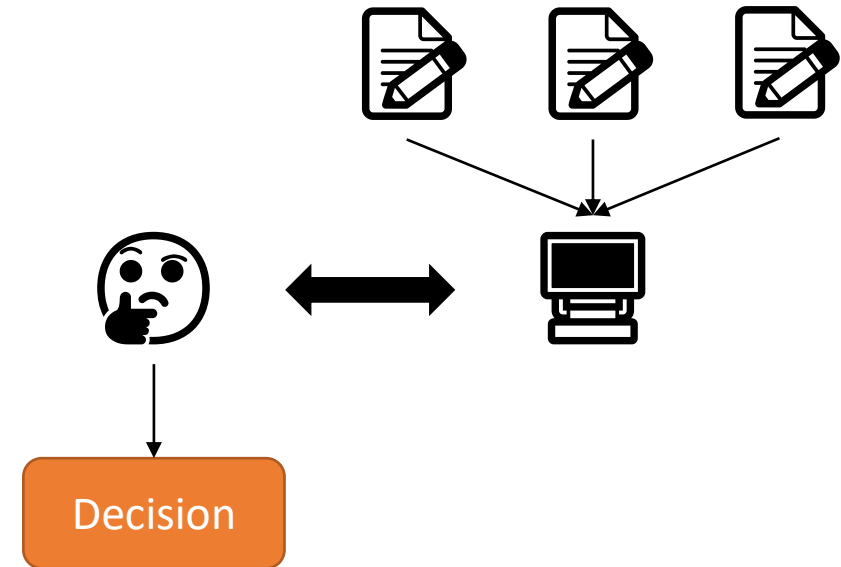
**Study 1:** Latent Agents
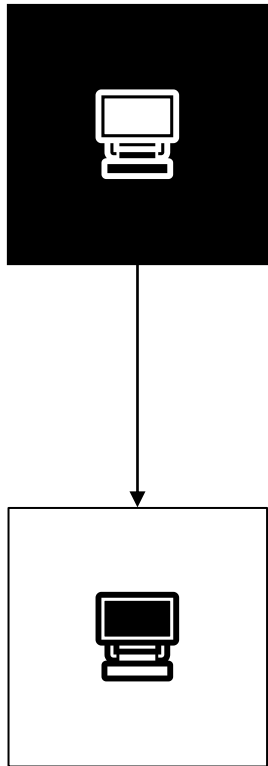
Study 2: Virtual Agents

Study 3: Embodied Agents

# Decision Support Systems (DSS)

- Used daily!

- Any AI system that assists the sensemaking process

- Exists in multiple domains: navigation, recommender systems, medicine diagnostics
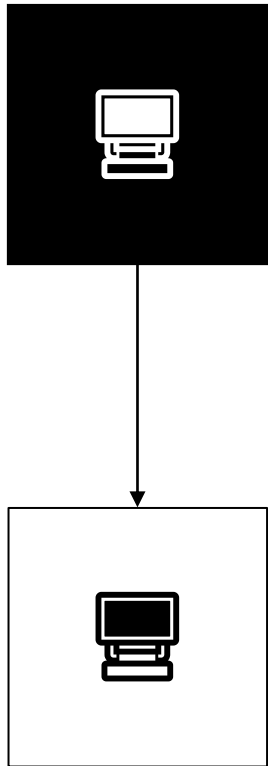
# From Black to White

- DSS are sophisticated and complex
  - Challenging even for experts to understand[a]
- Address with *transparency*[b] and *customizability*[c]
  - Increases trustworthiness and adoption
- New issue: Increasing and improperly calibrated trust

a V. Arnold, P. A. Collier, S. A. Leech, and S. G. Sutton. Impact of intelligent decision aids on expert and novice decision-makers judgments. Accounting & Finance, 44, 1. 2004.

b V. Arnold, N. Clark, P. A. Collier, S. A. Leech, and S. G. Sutton. The differential use and effect of knowledge-based system explanations in novice and expert judgment decisions. Mis Quarterly. 2006

c B. P. Knijnenburg, S. Bostandjiev, J. O'Donovan, and A. Kobsa. Inspectability and control in social recommenders. Conference on Recommender systems. 2012.

# From Black to White

- DSS research done: complacency, performance degradation, situation awareness
  - DSS feature X → metric Z
- Arnold et al. suggest knowledge plays a large factor in judgment and performance[a]
  - DSS feature X → knowledge → metric Z
- *Could white-box features harm knowledge?*

a V. Arnold, N. Clark, P. A. Collier, S. A. Leech, and S. G. Sutton. The differential use and effect of knowledge-based system explanations in novice and expert judgment decisions. Mis Quarterly. 2006
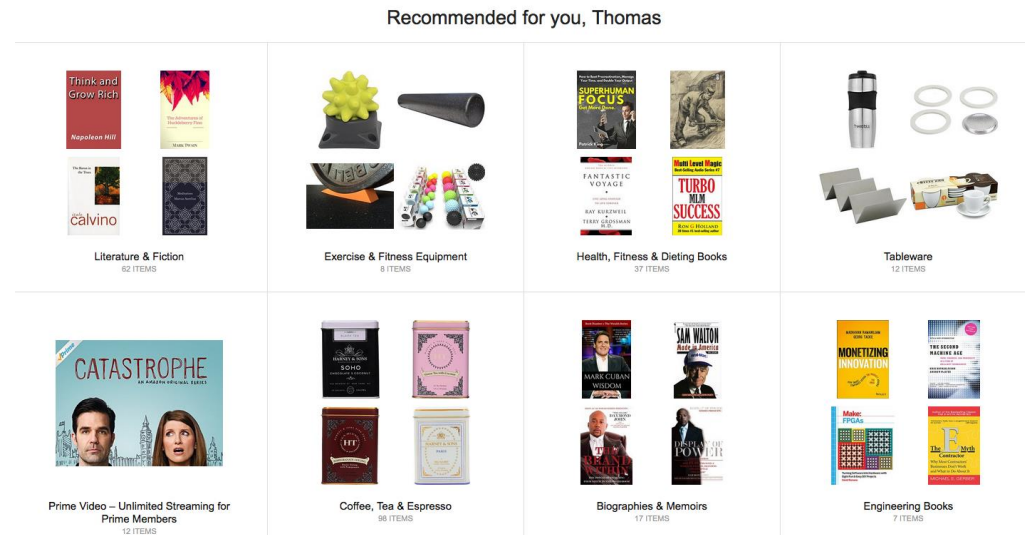
# Study 1: Latent Agents

## Research Question

How do latent agents affect domain knowledge acquisition and retention, and how is it affected by lowered reliability?

# Task Context

DSS tasks can be either subjective or objective



**Success = Satisfaction**
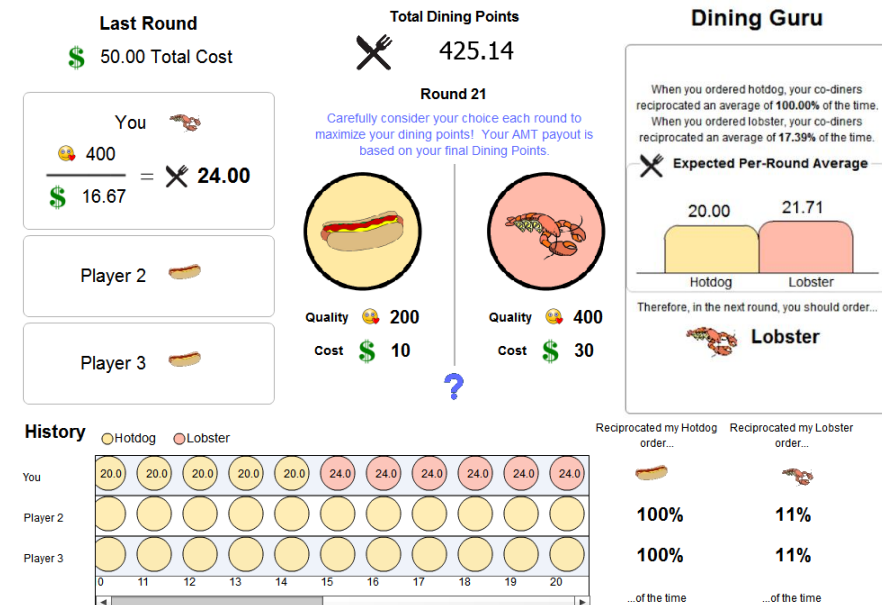


**Success = Performance**

# Task Selection: Recommender Systems

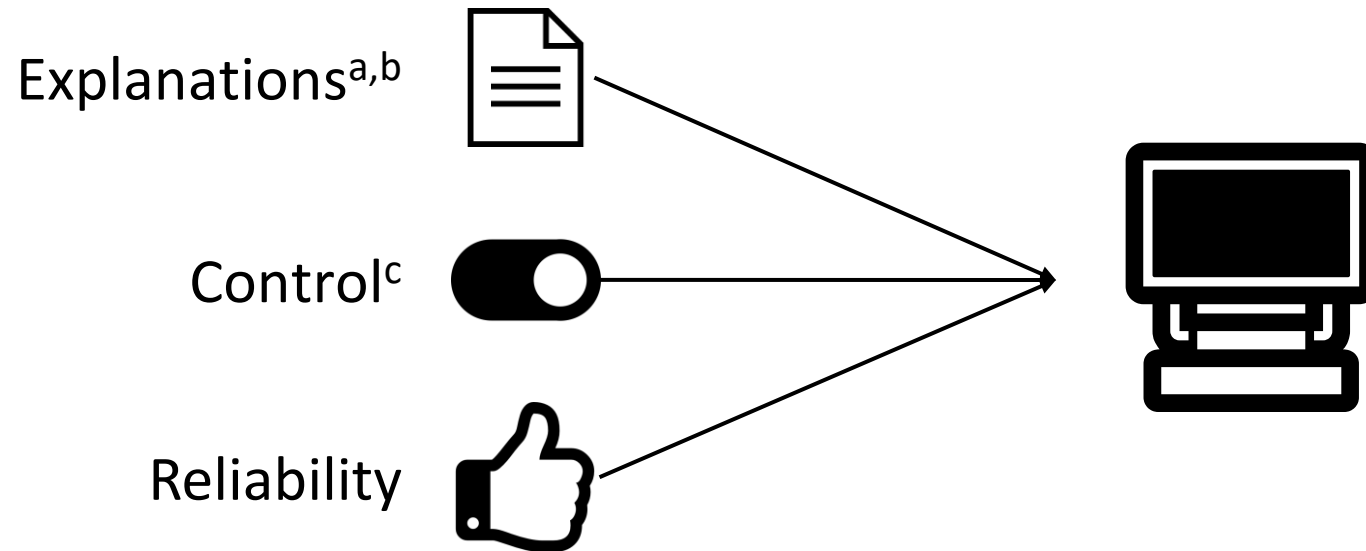DSS tasks can be either subjective or objective



**Success = Satisfaction**

*Movie Recommendation Task*

Make a watchlist from recommended movies

Predicted to forget knowledge about movies

**Success = Performance**

*Diner's Dilemma Game*

Eat the most while paying the least

Predicted to learn rules of the game

# Independent Variables



Explanations[a,b]

Control[c]

Reliability

These have granularities, but we'll simplify with dichotomies

a N. Tintarev and J. Masthoff. A survey of explanations in recommender systems. In Data Engineering Workshop, 2007 IEEE 23rd International Conference on.
b V. Arnold, N. Clark, P. A. Collier, S. A. Leech, and S. G. Sutton. The differential use and effect of knowledge-based system explanations in novice and expert judgment decisions. Mis Quarterly. 2006
c B. P. Knijnenburg, S. Bostandjiev, J. O'Donovan, and A. Kobsa. Inspectability and control in social recommenders. Conference on Recommender systems. 2012.

# Study Design
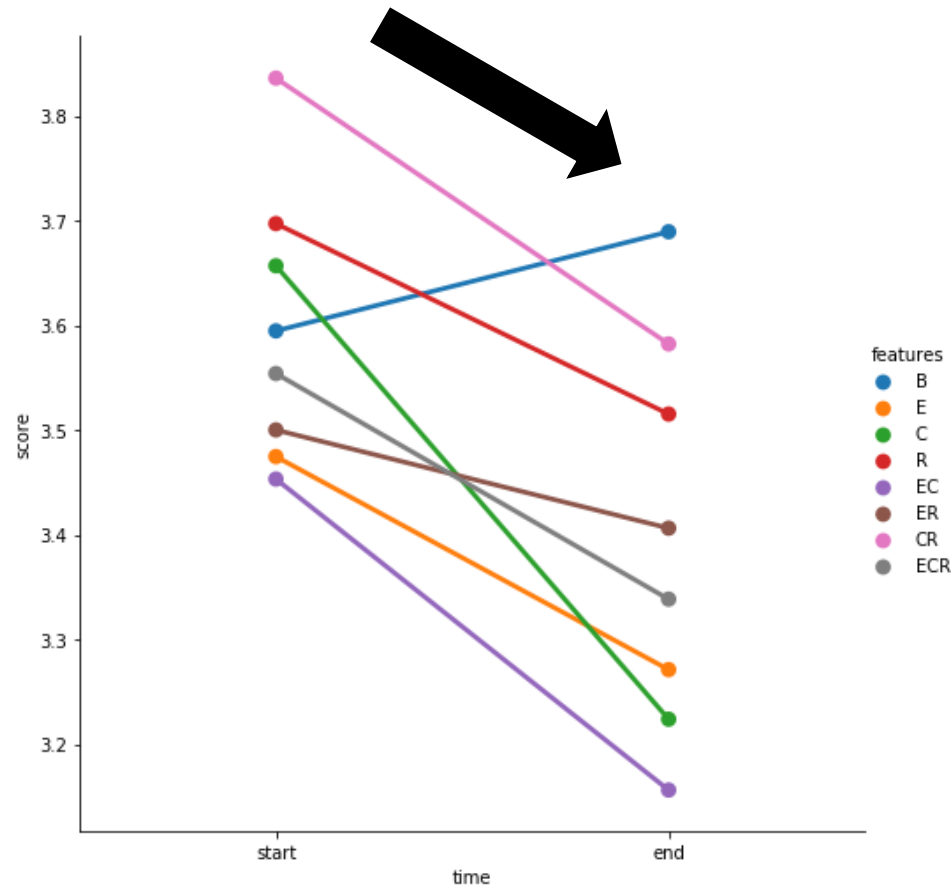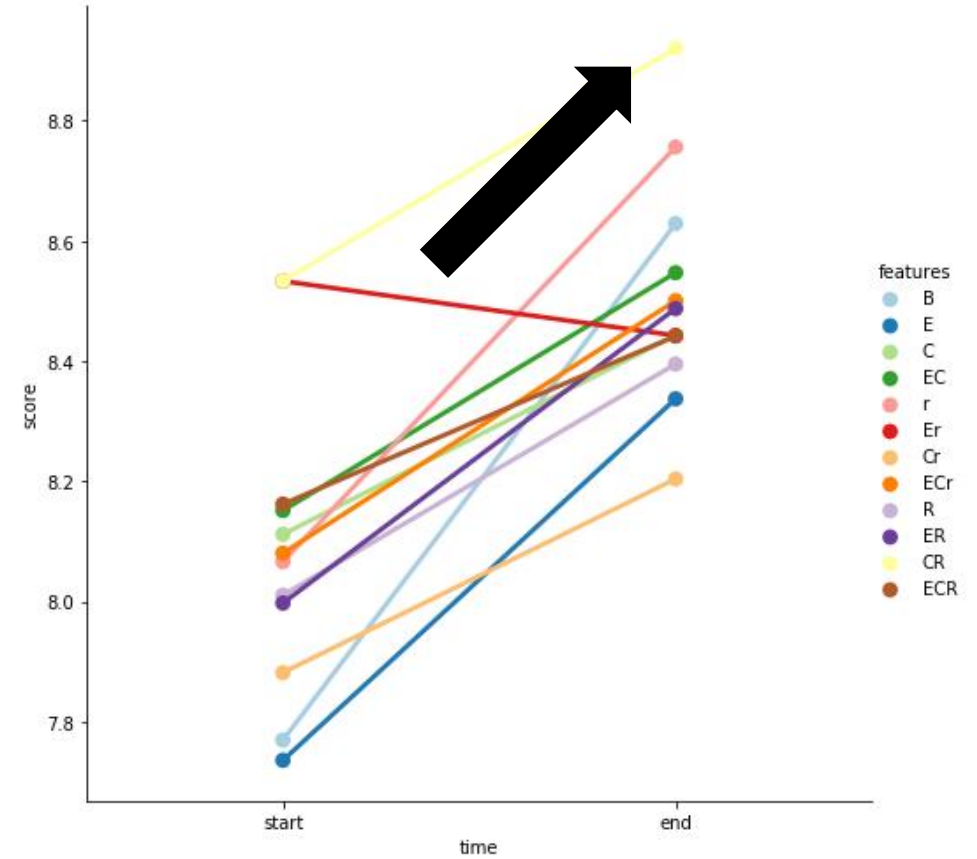
# Knowledge changes over time (expected)



Movie (loss)

Diner (gain)

# Potential gains are best with black boxes



Movie (loss)                    Diner (gain)

# Worst culprit: control settings

## Movie (n = 526)

|  | B | p-value |
|---|---|---|
| **E**xplanation | -0.298 | 0.088 |
| **C**ontrol | -0.527 | **0.002**** |
| **R**eliability | -0.391 | **0.022*** |
| **EC** | 0.549 | **0.026*** |
| **ER** | 0.501 | **0.044*** |
| **CR** | 0.571 | **0.019*** |
| **ECR** | -0.714 | **0.041*** |

## Diner (n = 529)

|  | B | p-value |
|---|---|---|
| **E**xplanation | -0.464 | **0.015*** |
| **C**ontrol | -0.563 | **0.004**** |
| **R**eliability | -0.477 | **0.025*** |
| **EC** | 0.568 | **0.037*** |
| **ER** | 0.327 | 0.27 |
| **CR** | 0.533 | 0.078 |
| **ECR** | -0.5 | 0.239 |

White box harms mitigated by low reliability
Task features led to different interactions; more research needed

# Summary



Feature-filled agents led to worsened learning – over-trust

Control was the worst culprit (corrupts the information retrieval process)

Lowering reliability hampers over-trust; users forced to exercise judgment

Caveat: lowered reliability itself can also cause wrong insights to be formed

Different tasks led to different feature effects (some overlap)

Task context remains relevant for interpreting outcomes

# Dissertation Comparative Analysis

Study 1: Latent Agents

**Study 2: Virtual Agents**

Study 3: Embodied Agents

# Physical agents in the world

Fire

Navigation

Reconnaissance

Domain e.g.: IoBT
(little human intervention)



Emergent collaboration with physical non-controllable entities
Hard to study in a safe manner

# Changes from latent agents

Supervisor

Agent

→

Teammates

Independent workers

- New considerations[a]
  - Non-supervisory role
  - Non-discrete decision making
  - Non-latent agents (virtuality simulating physicality)
- Scenarios difficult to study; simulation used for safety
- How does reliability, trust, and performance interact in simulated physical domains? Often not studied

**a** M. Endsley. From Here to Autonomy: Lessons Learned From Human-Automation Research. Human Factors, 59, 1. Feb 2017

# Study 2: Virtual Agents

## Research Questions

How is human performance affected by varying reliability in collaborating agents in a continuous pursuit task?

How do human individual differences mediate perception of agents, situation awareness, and performance in continuous domains?

# Task Context

- Select task that models agent capabilities and usage
  - ARL collaboration: aerial drones for multi-agent systems[a,b]
- Target recognition and pursuit
  - Not exclusive to IoBT
  - Fundamental capabilities for physical tasks[b]
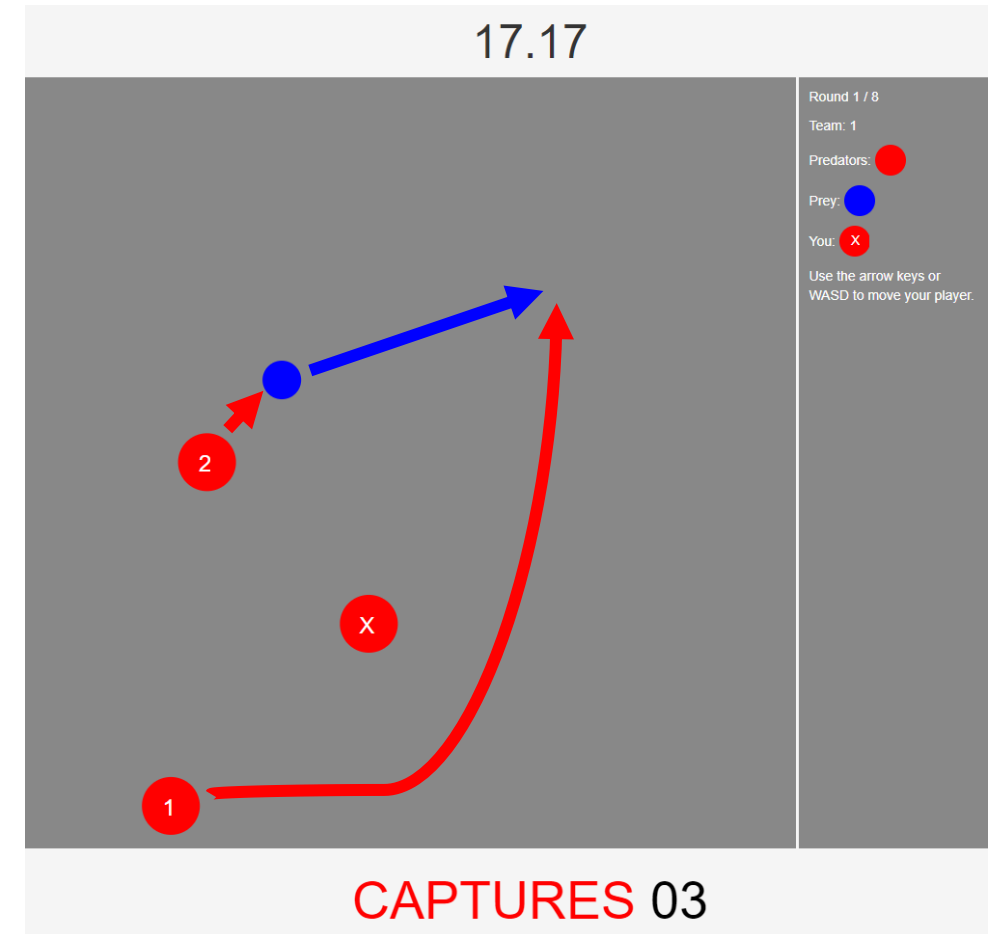  - Covers large part of multi-agent interaction and human-robot interaction

**a** D. Asher and S. Barton. 2018. Reinforcement learning framework for collaborative agents interacting with soldiers in dynamicmilitary contexts.
**b** D. Asher, E. Zaroukian, and S. Barton. 2018. Adapting the Predator-Prey Game Theoretic Environment to Army Tactical Edge Scenarios with Computational Multiagent Systems

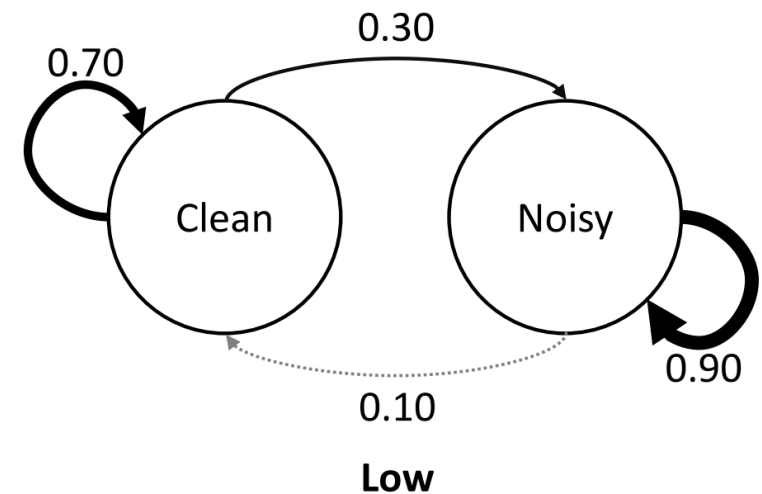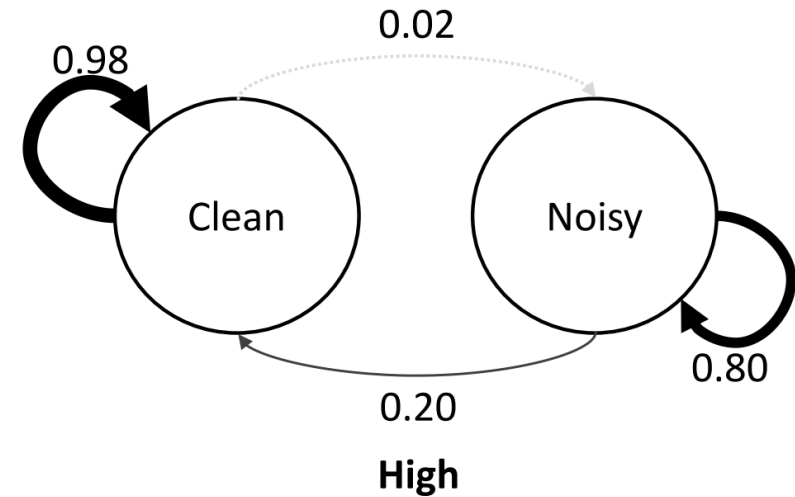# Task Selection: Predator-Prey Game

- Continuous pursuit task
  - Goal: capture the prey
- Asymmetrical and heterogeneous capabilities to engender collaboration
  - AI prey is faster
  - AI predators use different pursuit strategies (chase or intercept)
- Participants must be aware of their teammates for proper strategy-building
- Can situation awareness (SA) be affected by good performing agents?

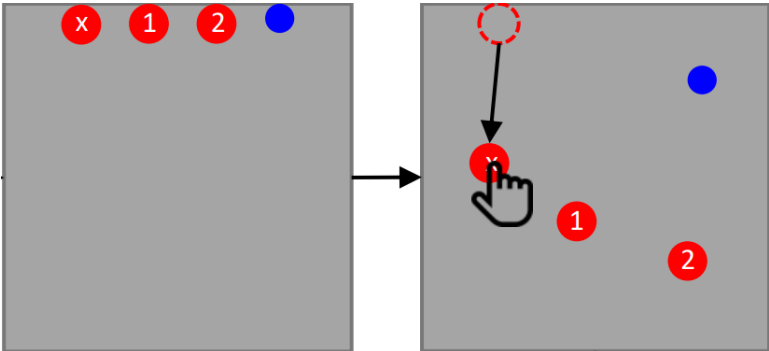# Independent Variable: Reliability

- How do we design reliability?
  - Goal: simulate environmental/system randomness
  - 2 states
    - Clean: no movement interference
    - Noisy: movement interference
  - 2 MDPs
    - High and Low
    - Both able to chase the prey; Low (i.e., "Good Enough") sometimes acts erratically
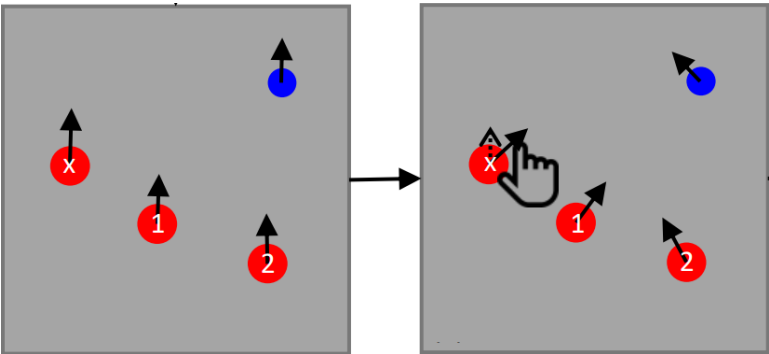
# Dependent Variables

### Situation Awareness (SA)
**State Reconstruction[f]**



Position

Direction

### Surveys
**Individual Differences (ID)/Subjective Perceptions (SP)**

| Factor | Construct | Survey |
|---|---|---|
| ID | **Motivation** | Intrinsic Motivation Inventory[a] |
| ID | **Trust Propensity** | Propensity to Trust Inventory[b] |
| ID | **Complacency Potential** | Automation-induced Complacency Potential Scale[c] |
| SP | **Trust in Agent** | Trust in Automated Systems[d] |
| SP | **Workload** | NASA Task Load Index[e] |

**a** Ryan and Deci, 1982 **b** Jessup et al., 2019 **c** Merritt et al., 2019 **d** Jian et al., 2000 **e** Hart and Staveland, 1988 **f** Endsley, 1988.

# Study Design

# Non-significant changes in SA (mostly)

| SA Probe | Base p-val | Change p-val |
|---|---|---|
| Position – Human | 0.81 | 0.41 |
| Position – Chaser | 0.84 | 0.89 |
| Position – Interceptor | 0.33 | 0.23 |
| Position – Prey | 0.54 | 0.20 |
| Direction – Human | 0.30 | 0.73 |
| Direction – Chaser | 0.47 | 0.29 |
| Direction – Interceptor | 0.69 | 0.35 |
| **Direction – Prey** | **0.049*** | 0.13 |



Direction Prey

Base p = 0.049*
Change p = 0.13

# Minimal changes in human performance

# Minimal changes in human performance

# Minimal changes in human performance



High

Low

Low reliability chaser improved performance
(task factors in play)

Agent
- Human
- Chaser
- Interceptor

score

Time

# Final SEM: Differing trust effects

# Summary



Reduced reliability led to improved team performance without affecting SA
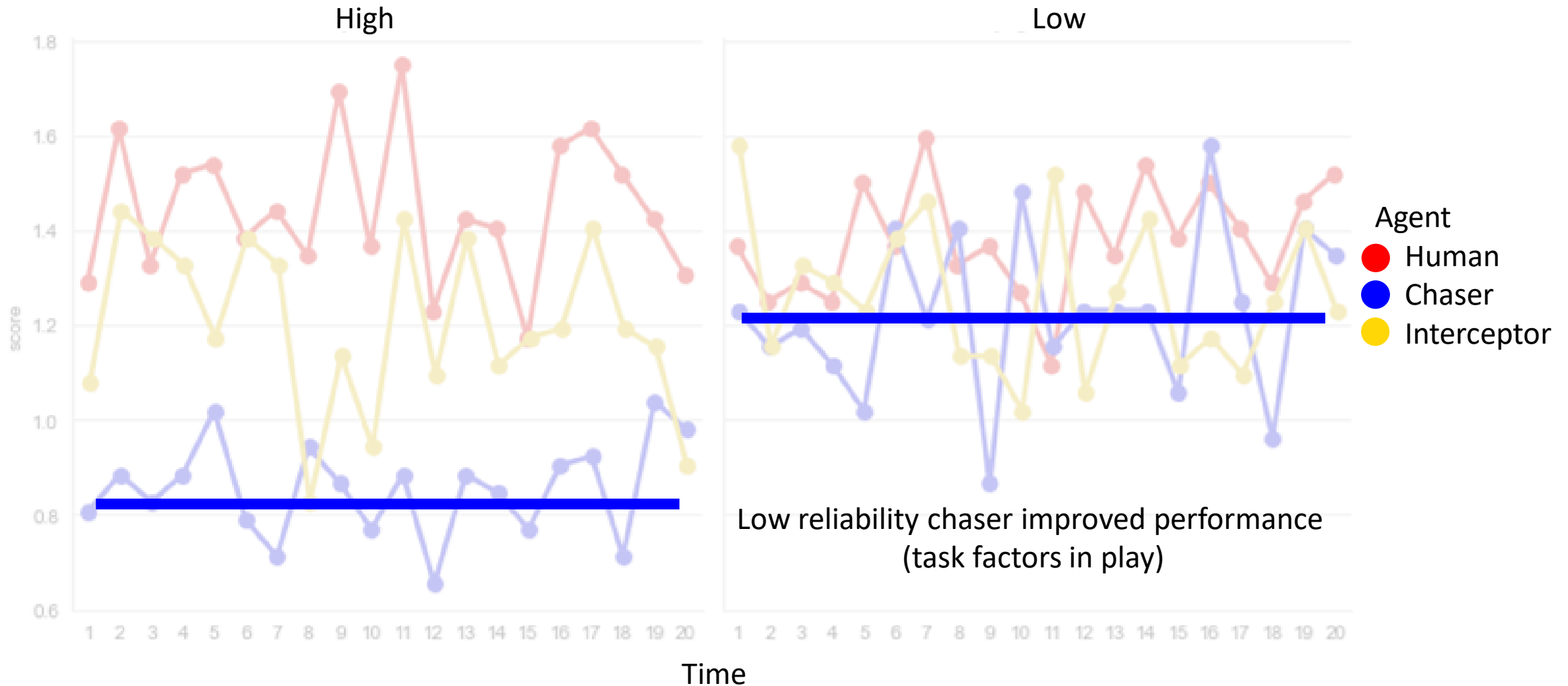
Performance may differ per domain: the PPG may have benefitted from erraticness

The resulting SEM matches models from HF research: physical tasks and continuous environments are subject to similar complacency dynamics

One exception: trust. Task type may have different influences

"Good enough" agents can be imperceptible, yet bring benefits

Could high reliability have minimal returns?

# Dissertation Comparative Analysis

**"Good Enough" Agents**

→ **Study 1**
Latent Agents

- "Good Enough" agents can help users exercise their judgment more often in influence tasks with decision support systems.
- The "gold" standard (explanations, control) could be "silver" instead.

→ **Study 2**
Virtual Agents

- "Good Enough" agents can contribute to improved situation awareness for behavioral predictions.
- Trust positively affected situation awareness; a new insight from non-supervisory control tasks.

→ **Study 3**
Embodied Agents

Study 1: Latent Agents

Study 2: Virtual Agents

**Study 3: Embodied Agents**

# Differences in virtuality vs. physicality

- Interaction with agents can be virtual or physical
- Physical interactions subject to social embodiment[a]
  - Human responses can be different[b]
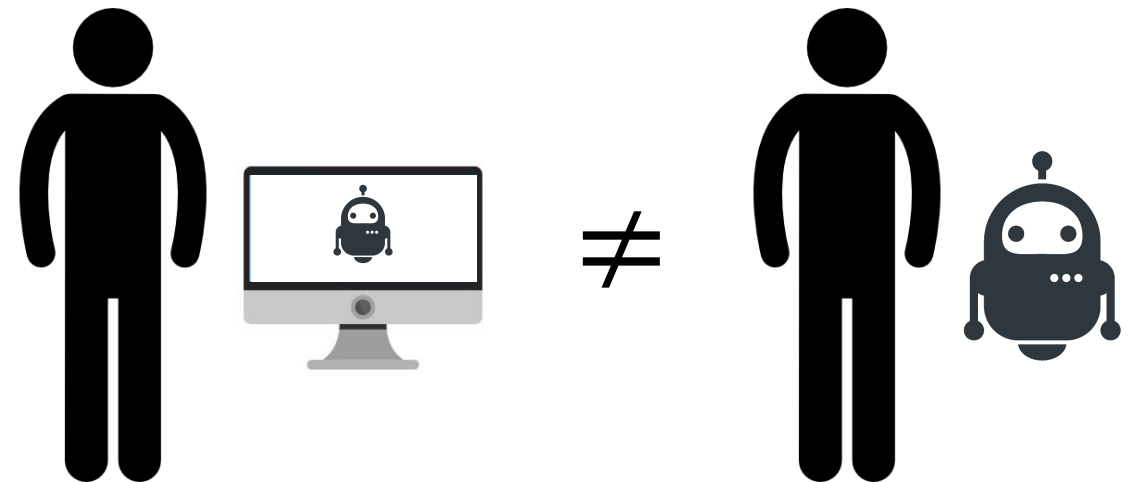- HAI research often researches tasks with virtual agents (see study 1 and 2)

$$\neq$$

a K. M. Lee. Presence, Explicated. Communication Theory, 14, 1. 2004. b L. W. Barsalou, P. M. Niedenthal, A. K. Barbey, and J. A. Ruppert. Social Embodiment. 2003.

# HRI and feasibility



- Physical agents and trust researched through human-robot interaction
  - Very limited: expensive and needs guarantees of safety

- Recent alternative: Virtual Reality[a]
  - Simulate any robot and related interactions

- Research on trust in VR upcoming

a O. Liu, D. Rakita, B. Mutlu, and M. Gleicher. Understanding Human-Robot Interaction in Virtual Reality. International Symposium on Robot and Human Interactive Communication. 2017.

# Embodiment and reliability interactions

- Can our previous reliability findings replicate with embodied agents?

- We can't really build or order a robot for this – use VR instead

# Study 3: Embodied Agents

## Research Questions

How does embodiment affect performance and trust towards an agent in a teamed physical task?

How does agent reliability interact with embodiment to affect human performance and perception?

# Task Context

- Common use cases for collaborative robotics (cobots)
  - Assembly, processing, packaging, un/loading, inspections
- Expectation: high reliability, human rarely in the loop (mostly automation)
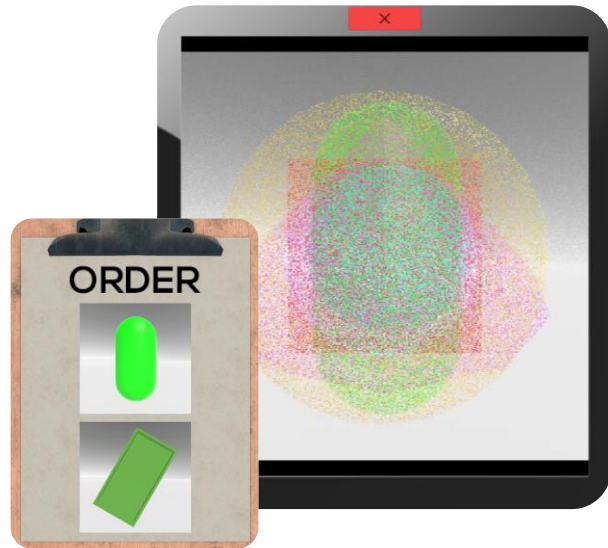- Introduce collaborative decision-making b/w humans and agents

# Task Selection: Warehouse

- Human and robot agent collaborate in signal detection theory-based game
  - Goal: Send out good packages, reject bad packages; as quick as possible
  - Robot gives a recommendation (package good or package bad – reliability varies within-subjects)
- 2 versions
  - Screen – telepresent agent
  - Virtual Reality – copresent agent
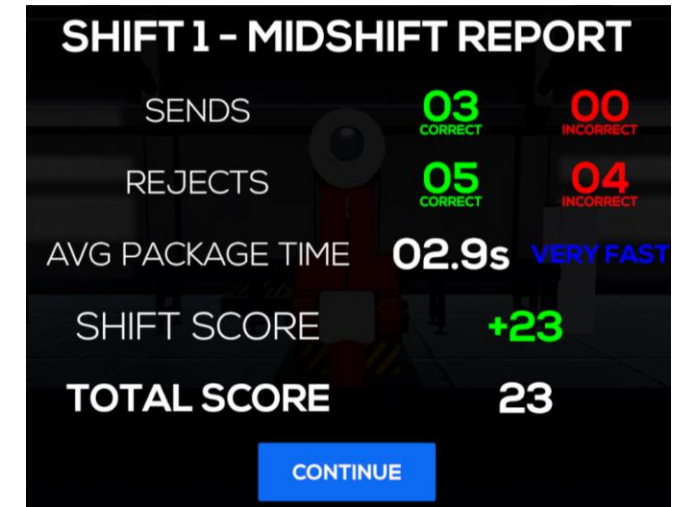
# Task Features Explored



**Signal Detection Theory (w/ time pressure)**
Given noisy feed, find target object
Adds uncertainty and reliance on robot recommendation
Asymmetric rewards



**Embodied Interactions**
VR participants could grab and interact with objects in the environment
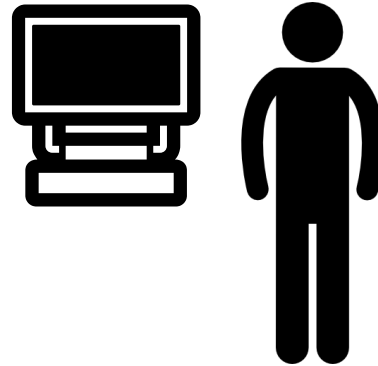Different levels of workload can incite distinct trust calibration



**Sporadic Feedback**
Participants would receive performance feedback at middle and end
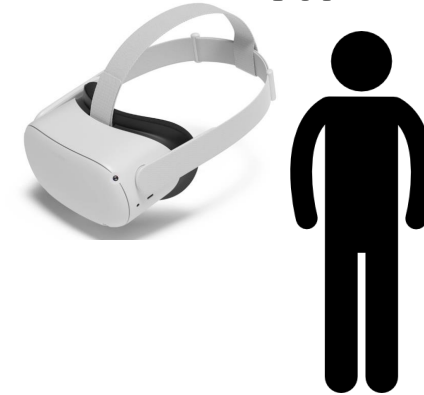Allows for discrete calibration points and time-based analyses

# Independent Variables

Representation
(2 levels)

**Screen**

**VR**

Reliability
(4 levels)

**Perfect
(100%)**

**Ideal
(91%)**

**Good
Enough
(75%)**

**No Info
(50%)**

# Measures

**Pre-Survey**         **Cognitive Reflection**, **Propensity to Trust**

**Task (Behavioral)**  **Accuracy**, **Reliance/Compliance**, Adherence, Switches, **Deferrals**
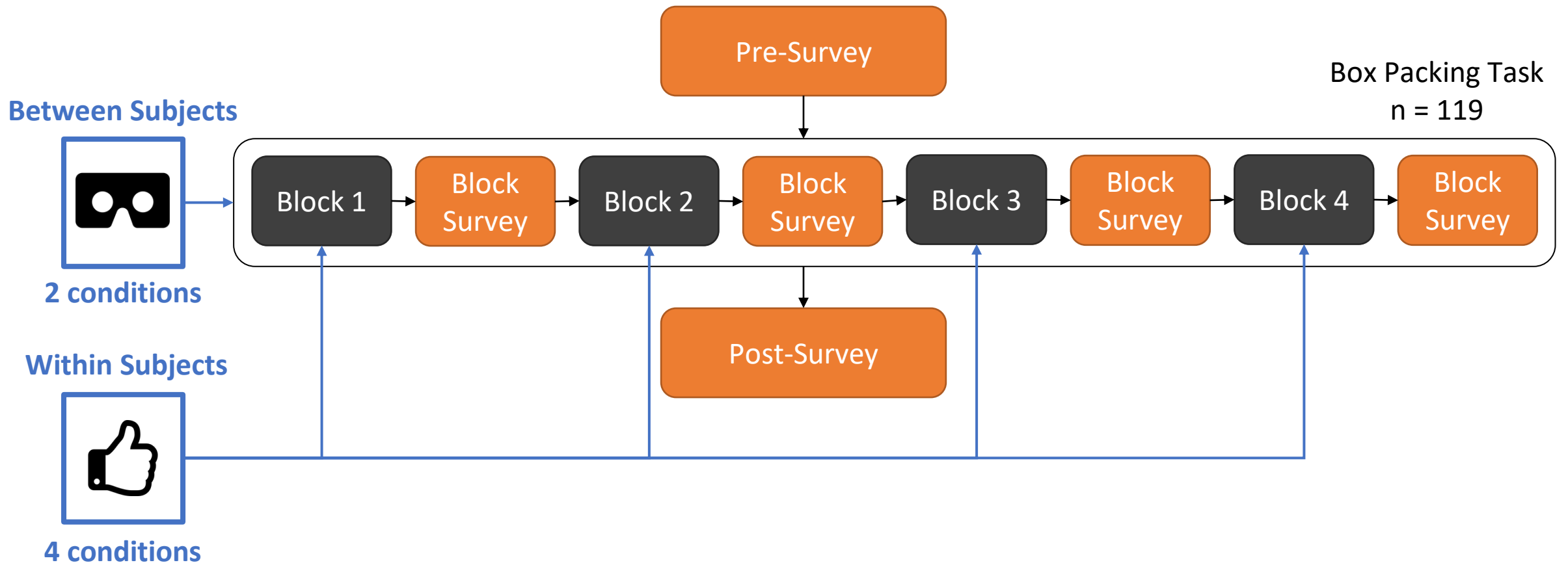
**Block Survey**       Workload, Information Preference, **Trust**, **Manipulation Check**

**Post-Survey**        Robot Ranking, Immersion

# Study Design

# Performance scales with reliability (expected)



Average Score per Participant

Performance expected to be a function of reliability (p < 0.001)

Best gains over time in **Good Enough** condition; not enough to outperform **Perfect** or **Ideal** (p < 0.001)

# Decision making at thresholds



Interesting sweetspots: **1 error** or **50% errors** (p < 0.001)
**Ideal:** Sends decreased
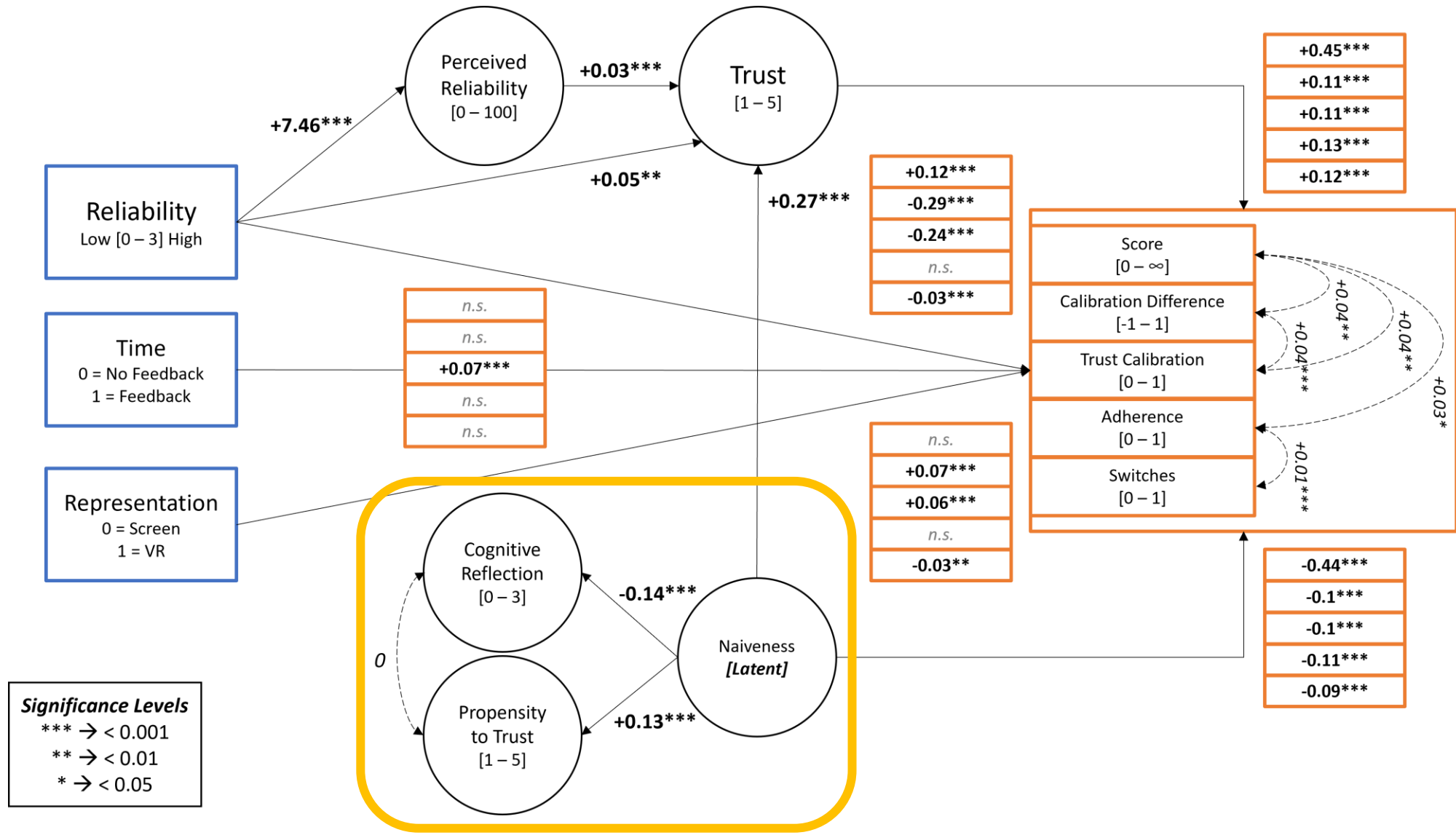**No Info:** Sends increased
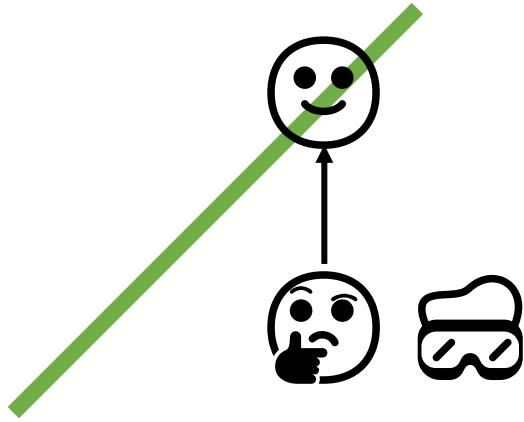
# Trust calibration improved with embodiment



VR had higher deferrals over time (p < 0.001),
**increasing overall trust calibration**

Reliability did not affect deferrals, thus it is a
property of embodied interactions

# Final SEM: Embodiment is good, "naiveness"



**Significance Levels**
*** → < 0.001
** → < 0.01
* → < 0.05

# Summary



Agent framing is important to establish a baseline of trust allocation

Embodiment allows for improved trust calibration following this allocation

Humans can determine reliability with little to no feedback, even in highly uncertain scenarios

This perceived reliability can then affect decision-making

Independent differences converged into a "naiveness" factor, worsening trust calibration and performance

"Naiveness" is mediated by trust, matching prior models establishing pre-dispositions and outcomes

# Dissertation Comparative Analysis

**"Good Enough" Agents**

**Study 1**
Latent Agents

- "Good Enough" agents can help users exercise their judgment more often in influence tasks with decision support systems.
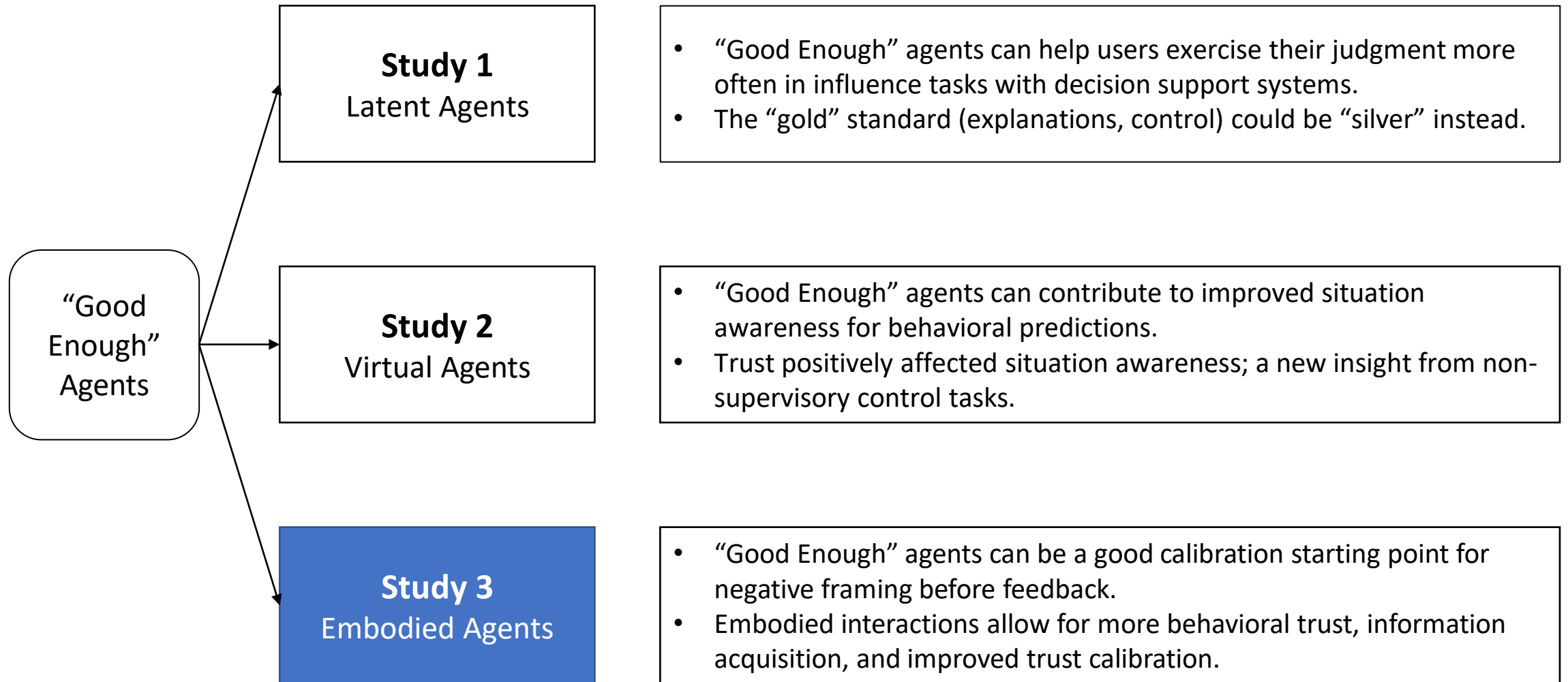- The "gold" standard (explanations, control) could be "silver" instead.

**Study 2**
Virtual Agents

- "Good Enough" agents can contribute to improved situation awareness for behavioral predictions.
- Trust positively affected situation awareness; a new insight from non-supervisory control tasks.

**Study 3**
Embodied Agents

- "Good Enough" agents can be a good calibration starting point for negative framing before feedback.
- Embodied interactions allow for more behavioral trust, information acquisition, and improved trust calibration.

# Closing Research Questions

*How does lowering agent reliability (i.e., "Good Enough" agents) affect task performance and human trust across different domains?*

Lower reliability *benefits*: prevent knowledge loss, exploit environmental interactions, and improve trust calibration under framing.

**These *benefits* can also describe performance and trust.**

# Closing Research Questions

*What is the relationship between human individual differences, agent reliability, and technology use?*

Individual differences **interplay** with SA, motivation, and workload.

These constructs affect outcomes on a domain-by-domain basis.

**Categorizing domains** may be a viable research approach.

# Closing Research Questions

*How do different task requirements interact with the human-AI interaction dynamic?*
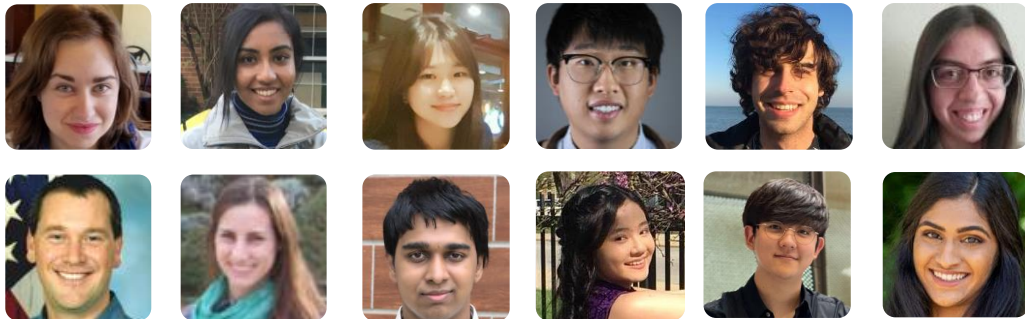
Tasks may require **time pressure, different invocation strategies, agent framing, agent representation, reliability thresholds**, with potential others.

Requirements can change how performance and trust develops.

**We should consider *not* looking at IV/DVs in a vacuum (when possible).**