

# Trust in and dependence on **imperfect** automation

Jessie Yang, PhD

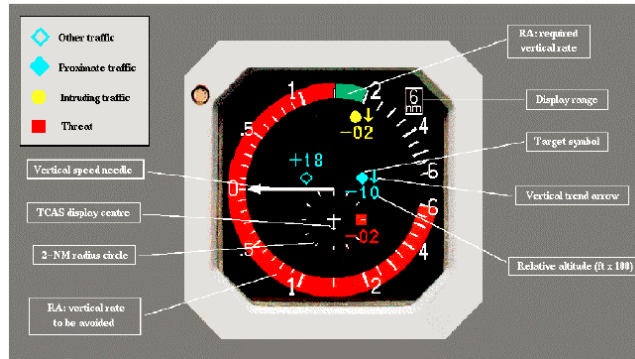
Industrial and Operations Engineering

University of Michigan, Ann Arbor

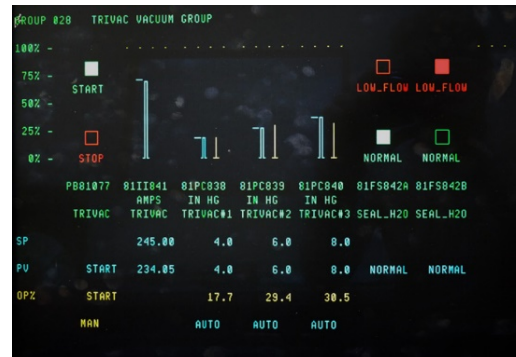


# Automation dependence and performance | Motivation

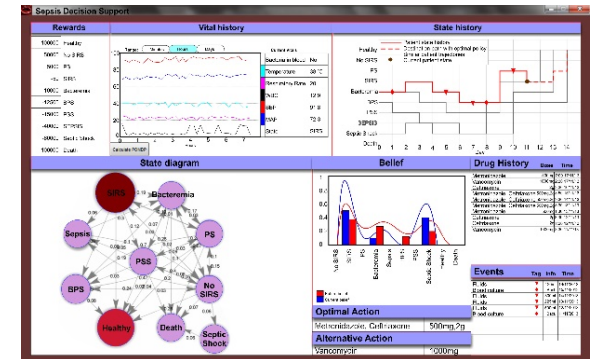
- Use of automation to assist human performance



Aircraft collision avoidance system



Alarm management



Clinical decision support system

- Ideally – performance gain; In reality, not always..



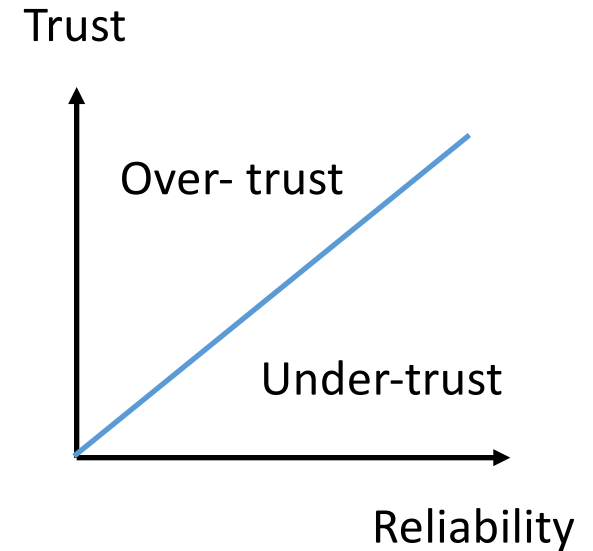
Grounding of the Royal Majesty, 1995



Crash of Korean air flight 801, 2001

# Automation dependence | Introduction

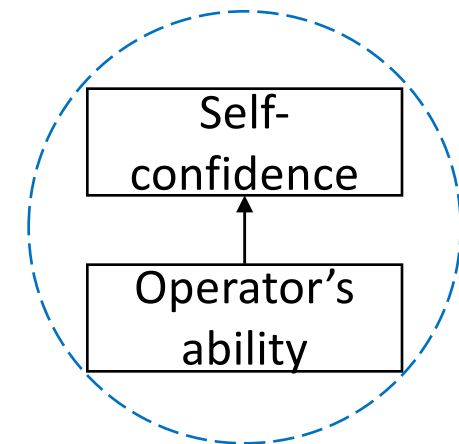
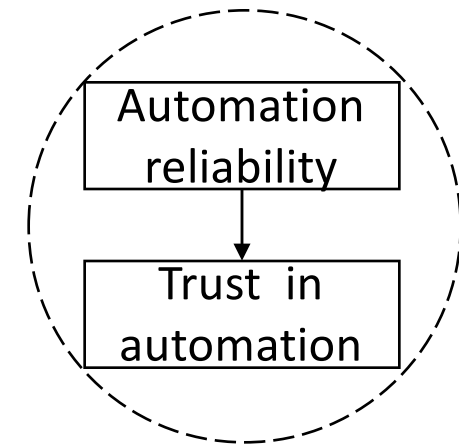
- Inappropriate use of automation<sup>1,2</sup>
- Two major reasons
  - Automation is sometimes **imperfect**<sup>3,5</sup>
  - Trust-reliability **miscalibration**<sup>4,5</sup>
- Trust in automation - belief, intention, **attitude**<sup>2</sup>, behavior
- Trust is an attitude, usage/dependence is a behavior



# Automation dependence | Two types of miscalibrations

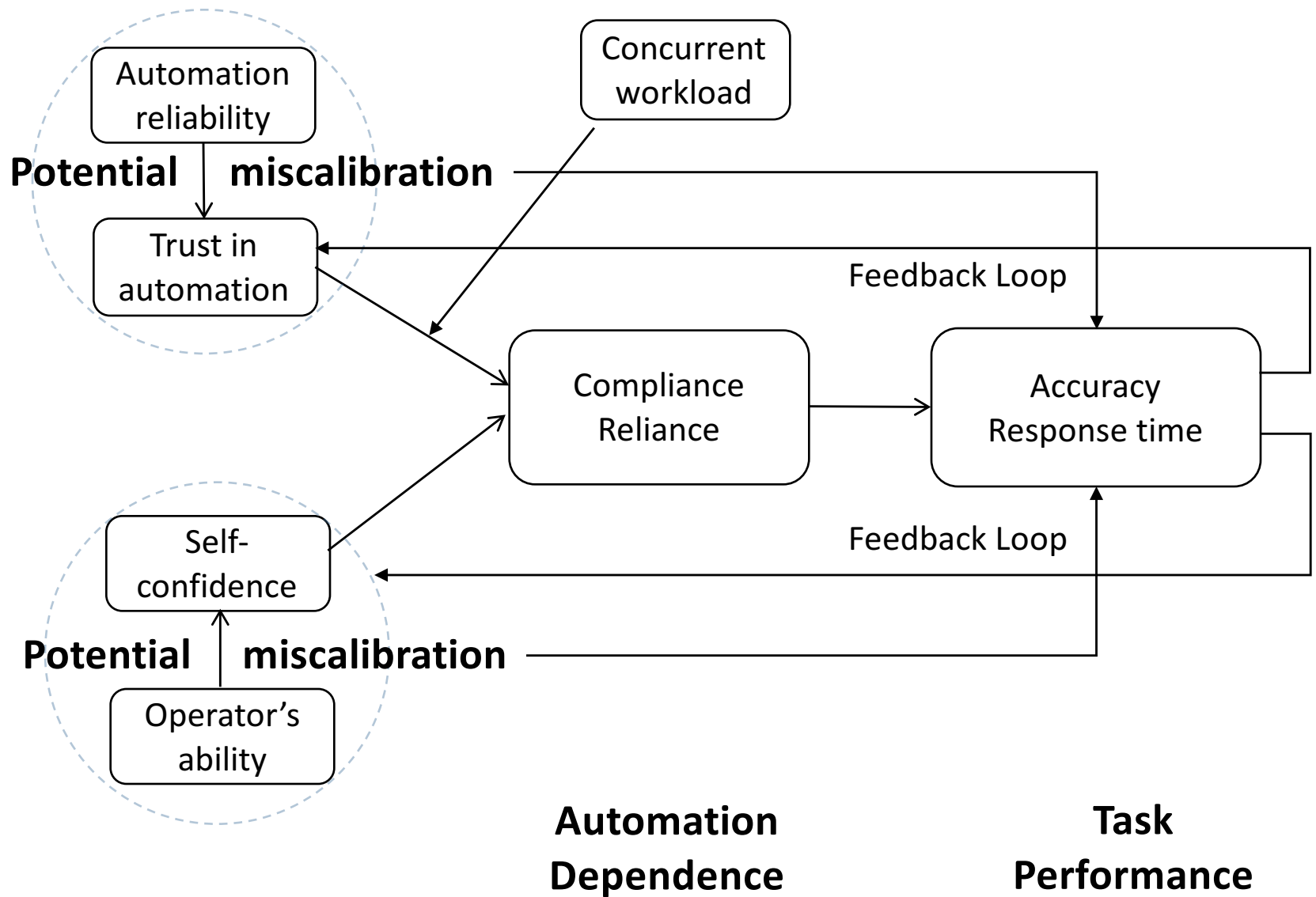
- Little attention on the human operator's ability
- Another factor: self-confidence in performing a task manually<sup>1,2,3</sup>
- Overconfidence is highly likely<sup>4</sup>, especially when tasks are difficult<sup>5</sup>
- Both types of miscalibrations should be modelled in human-automation interaction

Trust-reliability miscalibration

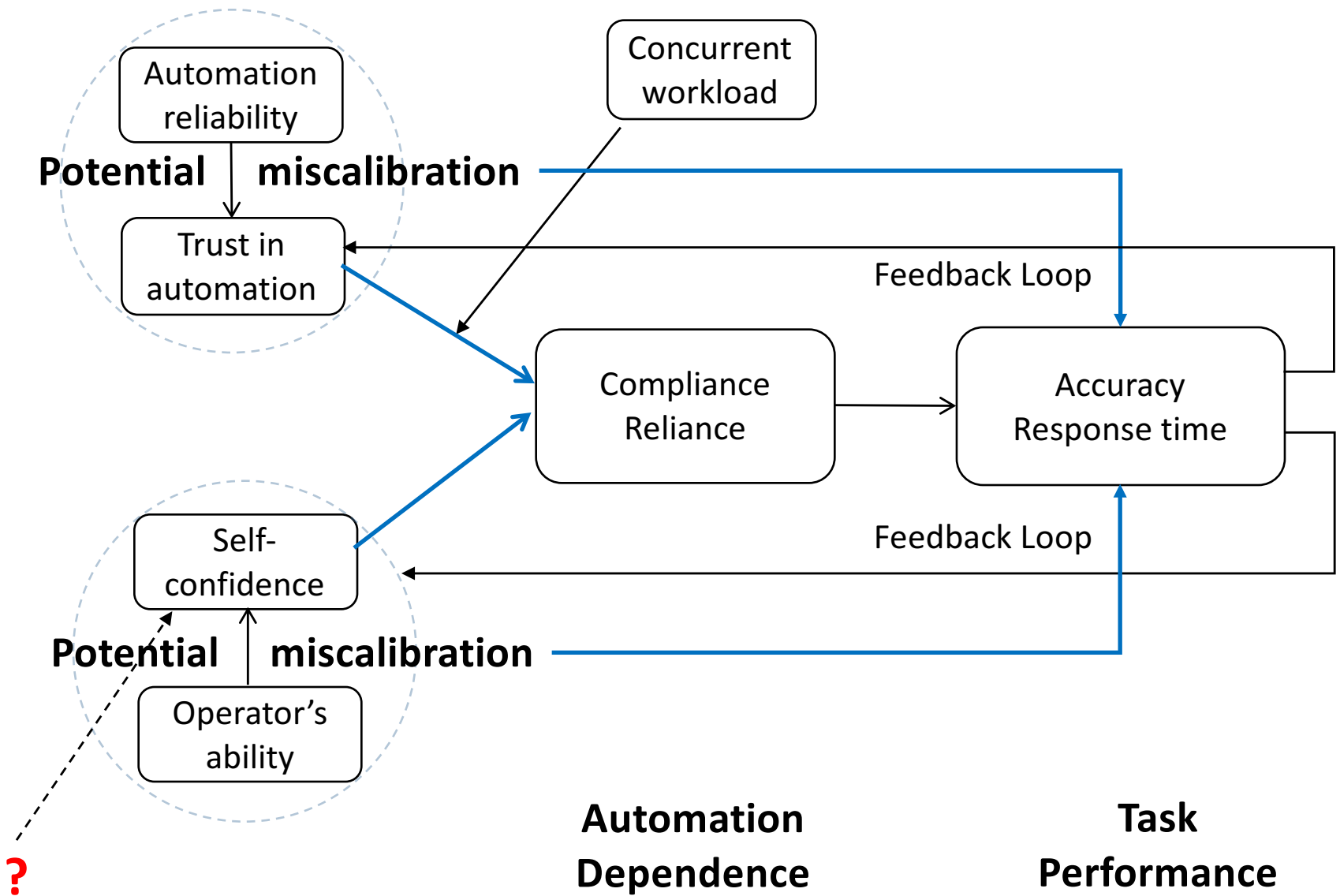


Confidence-ability miscalibration

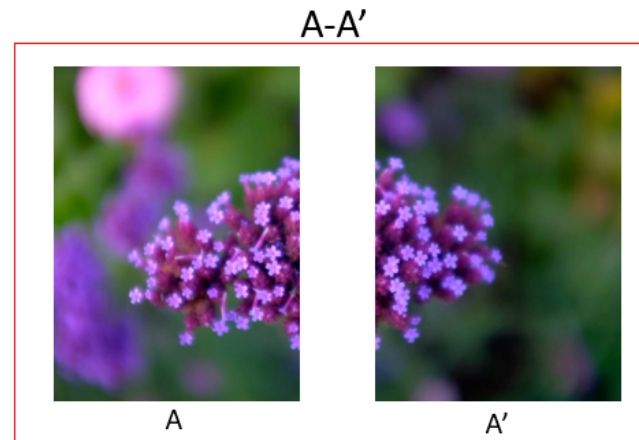
# Automation dependence | Research model



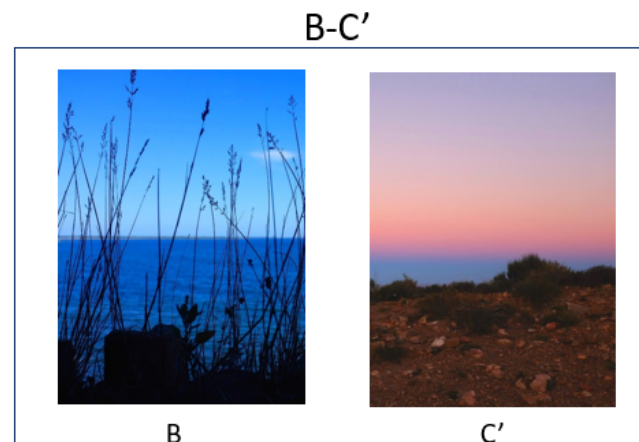
# Study 1 | Hypotheses



# Study 1 | Confidence-accuracy inversion<sup>1</sup>



Accuracy ↑  
Confidence ↓  
**Lower**  
overconfidence

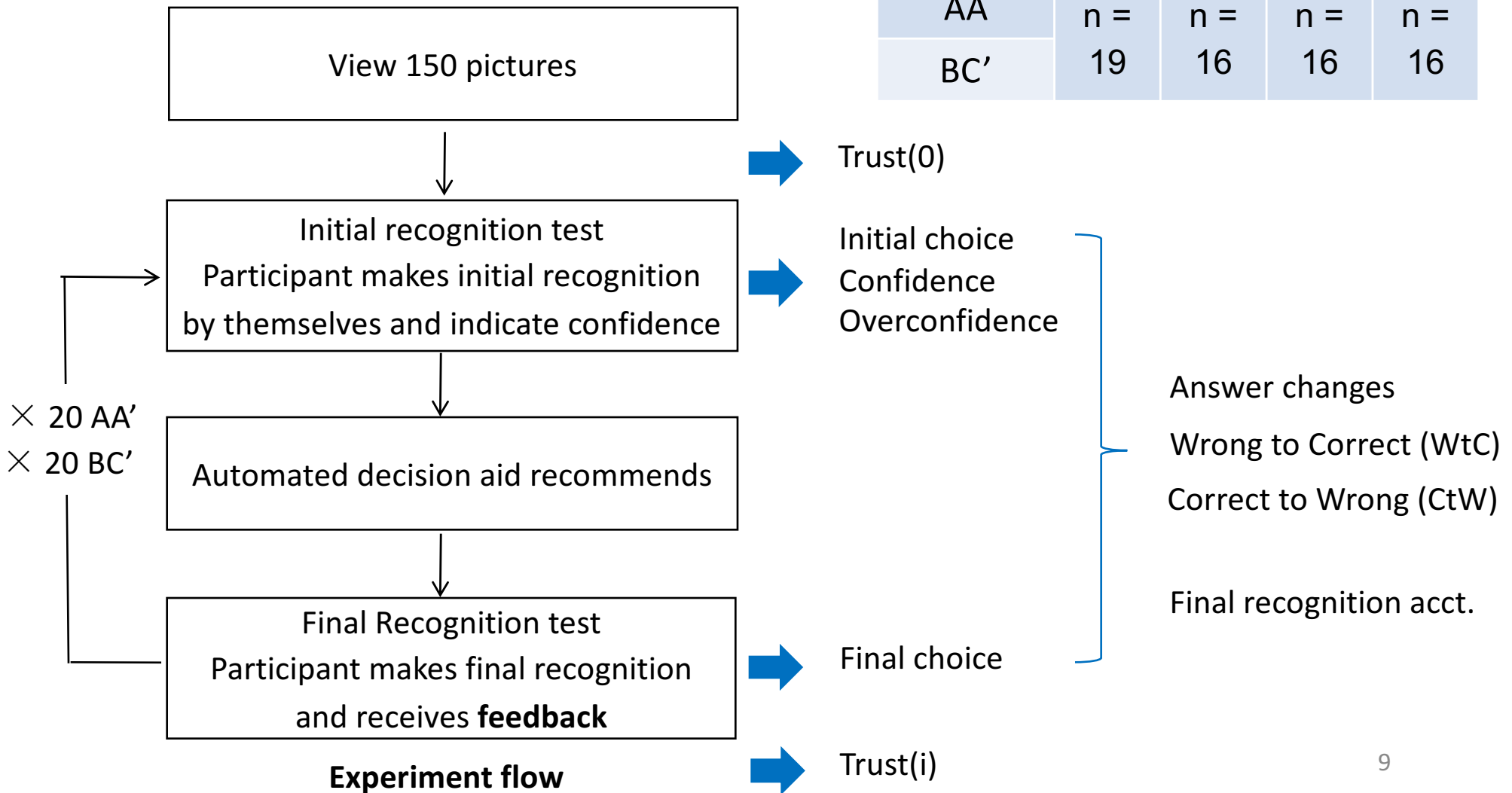


Accuracy ↓  
Confidence ↑  
**Higher**  
overconfidence

# Study 1 | Experimental Design

67 participants, age = 25.1 (SD = 3.8)

Reliability Picture	60%	70%	80%	90%
AA	n =	n =	n =	n =
BC'	19	16	16	16





## Study 1 | Results summary

**No. of answer change = 17.43 + 0.025 Trust<sup>\*\*\*</sup> - 0.045 Confidence<sup>\*\*\*</sup>**

**No. of answer change WtC = 22.81 + 0.038 Trust<sup>\*\*\*</sup> - 0.046 Confidence<sup>\*\*\*</sup>**

**No. of answer change CtW = 22.81 - 0.003 Trust - 0.043 Confidence<sup>\*\*\*</sup>**

- Self-confidence: strong and stable predictor of automation dependence

## Study 1 | Results summary

**No. of answer change** = 17.43 + 0.025 **Trust**<sup>\*\*\*</sup> – 0.045 **Confidence**<sup>\*\*\*</sup>

**No. of answer change WtC** = 22.81 + 0.038 **Trust**<sup>\*\*\*</sup> – 0.046 **Confidence**<sup>\*\*\*</sup>

**No. of answer change CtW** = 22.81 – 0.003 **Trust** – 0.043 **Confidence**<sup>\*\*\*</sup>

- Self-confidence: strong and stable predictor of automation dependence

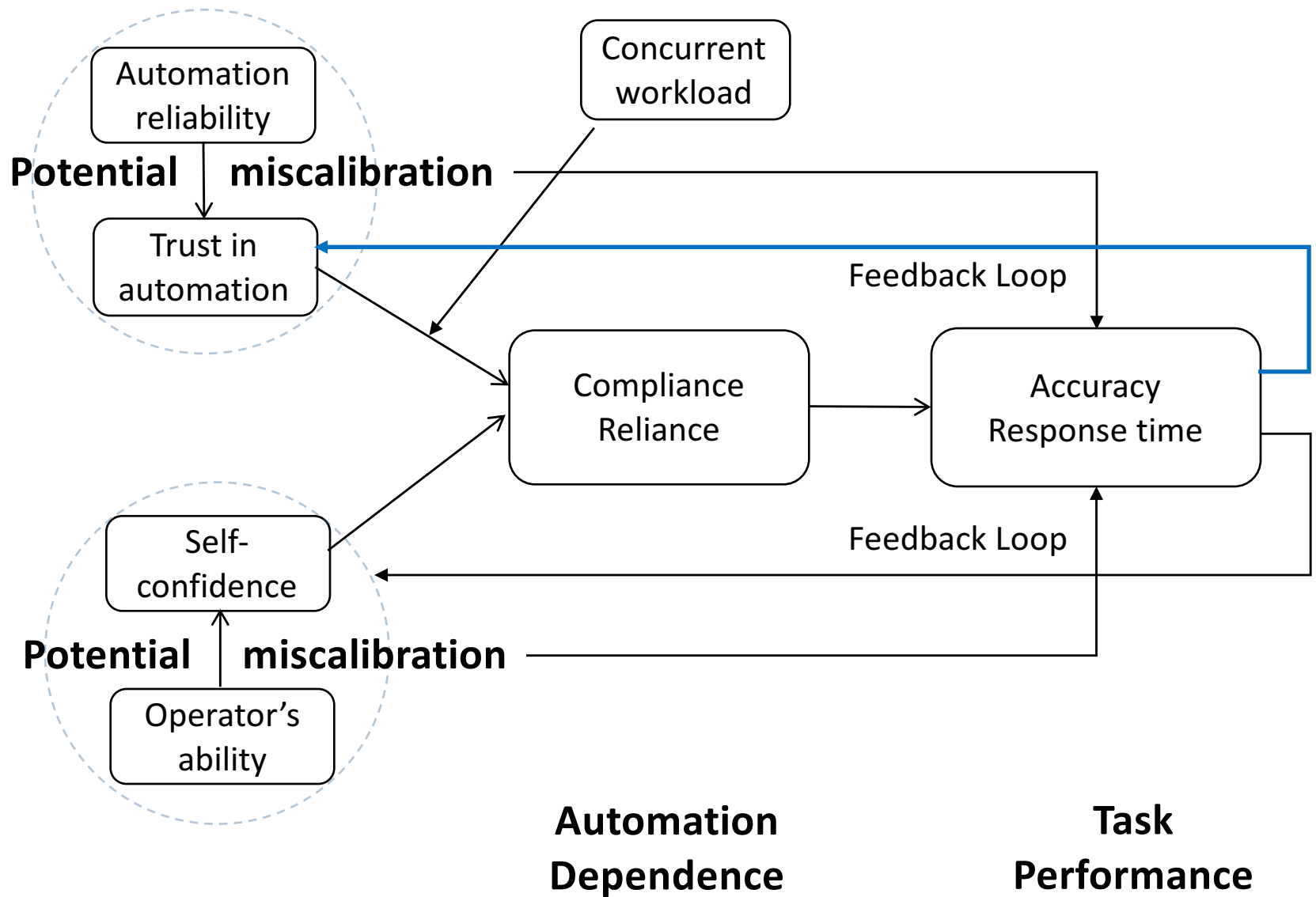
**Final recognition accuracy** = 19.94 + 0.56 **Initial recognition accuracy**<sup>\*\*\*</sup>

– [if reliability < ability] \* 0.38 **Overtrust**<sup>\*\*</sup>

– [if reliability > ability] \* 0.26 **Overconfidence**<sup>\*</sup>

- Both types of miscalibrations harms performance

# Automation dependence | Study 2



## Automation dependence | Study 2

Consider the following hypothetical situation



Mark



Clive

Mark and Clive had exactly the same medical condition. They were presented to the hospital with stomachache

## Consider the following hypothetical situation

Medical  
Decision Aid



Mark

Low risk

Endoscopy unnecessary



Clive

Low risk

Endoscopy unnecessary

Considering family history, age, dietary etc, a clinical decision support system suggested both Mark and Clive were at low risk of getting stomach cancer – endoscopy was unnecessary

# Consider the following hypothetical situation

Medical  
Decision Aid



Low risk  
Endoscopy unnecessary



pre-cancerous  
polyps removed

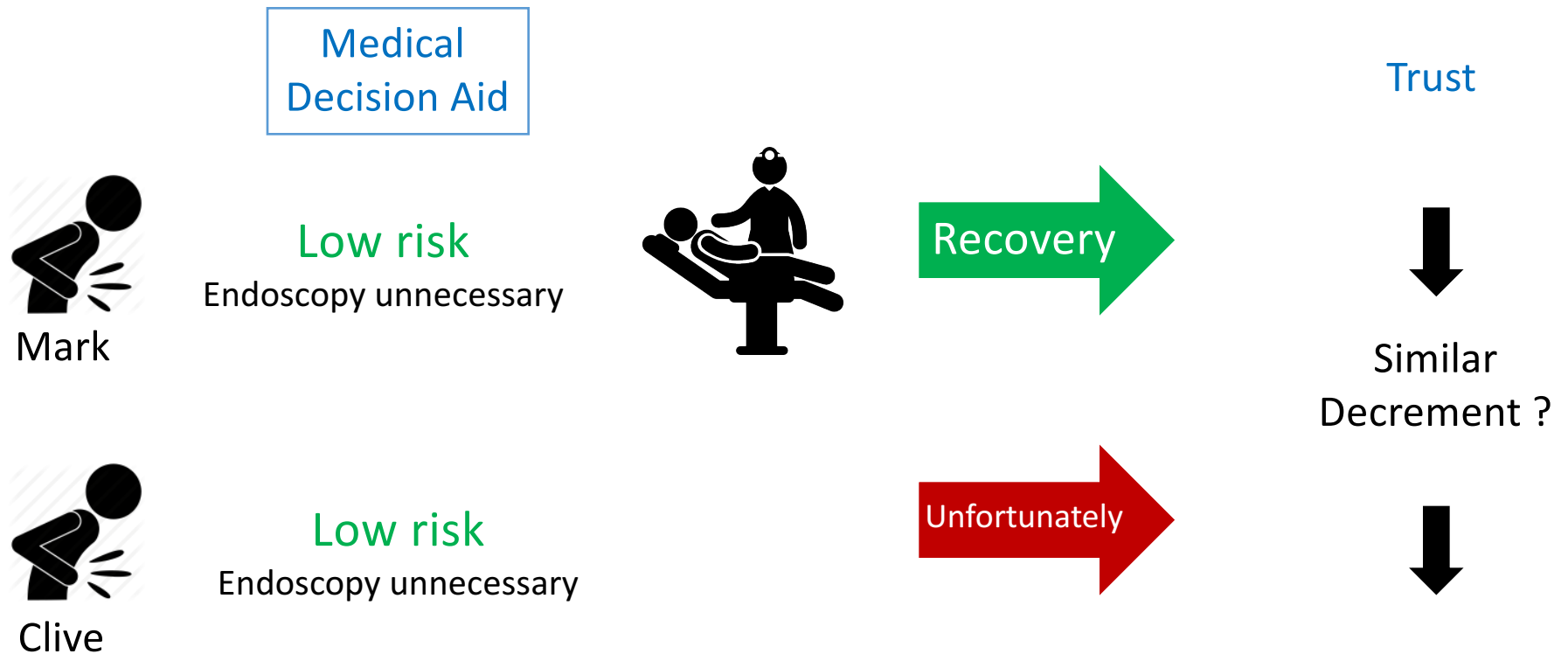


Low risk  
Endoscopy unnecessary

diagnosed with  
stomach cancer  
years later



# Consider the following hypothetical situation



In both cases, the medical decision aid made mistakes – we expect a decrement in trust. But **similar decrement?**

## Study 2 | Results

Cond.	Initial recognition	Recommendation	Final recognition	Trust adjustment
0: 000	Wrong	Wrong	Wrong	-4.0(0.5)
1: 001	Wrong	Wrong	Correct	
2: 010	Wrong	Correct	Wrong	2.0 (0.3)
3: 011	Wrong	Correct	Correct	2.7 (0.5)
4: 100	Correct	Wrong	Wrong	-6.2(1.2)
5: 101	Correct	Wrong	Correct	-4.3(0.6)
6: 110	Correct	Correct	wrong	
7: 111	Correct	Correct	Correct	1.5 (0.2)

Automation success → trust increases

Automation failure → trust decreases

\*\*\*  $p < .001$



## Study 2 | Results

Cond.	Initial recognition	Recommendation	Final recognition	Trust adjustment
<b>0: 000</b>	Wrong	<b>Wrong</b>	Wrong	-4.0(0.5)
1: 001	Wrong	Wrong	Correct	
<b>2: 010</b>	Wrong	<b>Correct</b>	Wrong	2.0 (0.3)
3: 011	Wrong	Correct	Correct	2.7 (0.5)
4: 100	Correct	Wrong	Wrong	-6.2(1.2)
<b>5: 101</b>	Correct	<b>Wrong</b>	Correct	-4.3(0.6)
6: 110	Correct	Correct	wrong	
<b>7: 111</b>	Correct	<b>Correct</b>	Correct	1.5 (0.2)

|Trust decrement| > |Trust increment|

\*\*\*  $p < .001$

## Study 2 | Results

Cond.	Initial recognition	Recommendation	Final recognition	Trust adjustment
0: 000	Wrong	Wrong	Wrong	-4.0(0.5)
1: 001	Wrong	Wrong	Correct	
2: 010	Wrong	Correct	Wrong	2.0 (0.3)
<b>3: 011</b>	<b>Wrong</b>	Correct	Correct	<b>2.7 (0.5)</b>
4: 100	Correct	Wrong	Wrong	-6.2(1.2)
5: 101	Correct	Wrong	Correct	-4.3(0.6)
6: 110	Correct	Correct	wrong	
<b>7: 111</b>	<b>Correct</b>	Correct	Correct	<b>1.5 (0.2)</b>

Automation successes lead to **greater** increment of trust, if a user is **less capable** of completing the task on his or her own. \*  $p < .05$

## Study 2 | Results

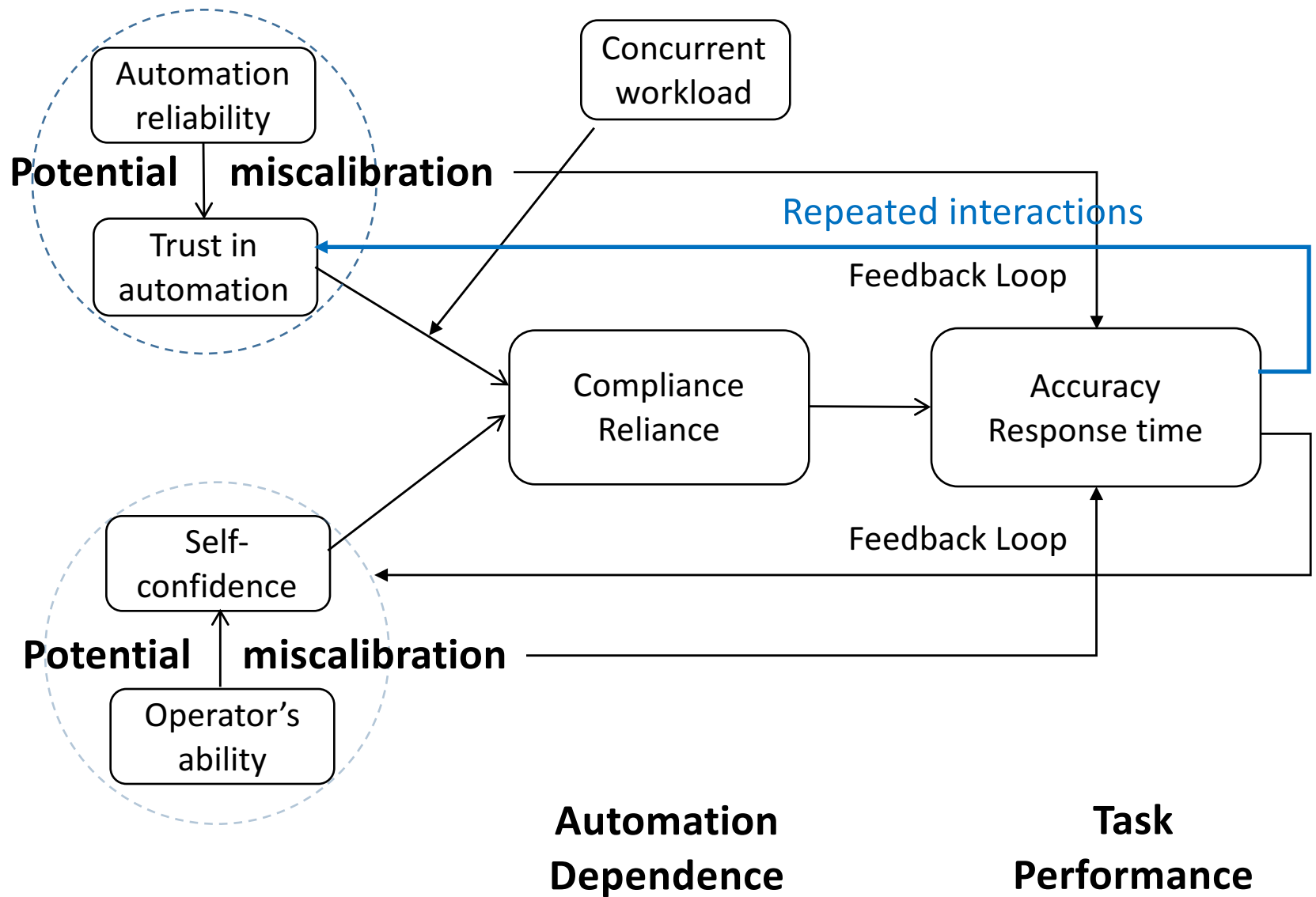
Cond.	Initial recognition	Recommendation	Final recognition	Trust adjustment
0: 000	Wrong	Wrong	Wrong	-4.0(0.5)
1: 001	Wrong	Wrong	Correct	
2: 010	Wrong	Correct	Wrong	2.0 (0.3)
3: 011	Wrong	Correct	Correct	2.7 (0.5)
<b>4: 100</b>	Correct	Wrong	<b>Wrong</b>	-6.2(1.2)
<b>5: 101</b>	Correct	Wrong	<b>Correct</b>	-4.3(0.6)
6: 110	Correct	Correct	wrong	
7: 111	Correct	Correct	Correct	1.5 (0.2)

Automation failures lead to **less** decrement of trust, if the outcome is not harmed. \*\*\*  $p < .001$

## Study 2 | Results

- $|\text{Trust decrement}| > |\text{Trust increment}|$
- Trust assessment is **not entirely rational**
  - **Not** benchmarked strictly against predetermined objective criteria
  - Contrast effect - based on **one's ability**: a correct recommendation is appreciated more if one cannot perform the task
  - Hindsight bias - based on **task outcomes**: a wrong recommendation is “forgiven” if it does not harm

# Automation dependence | Study 3

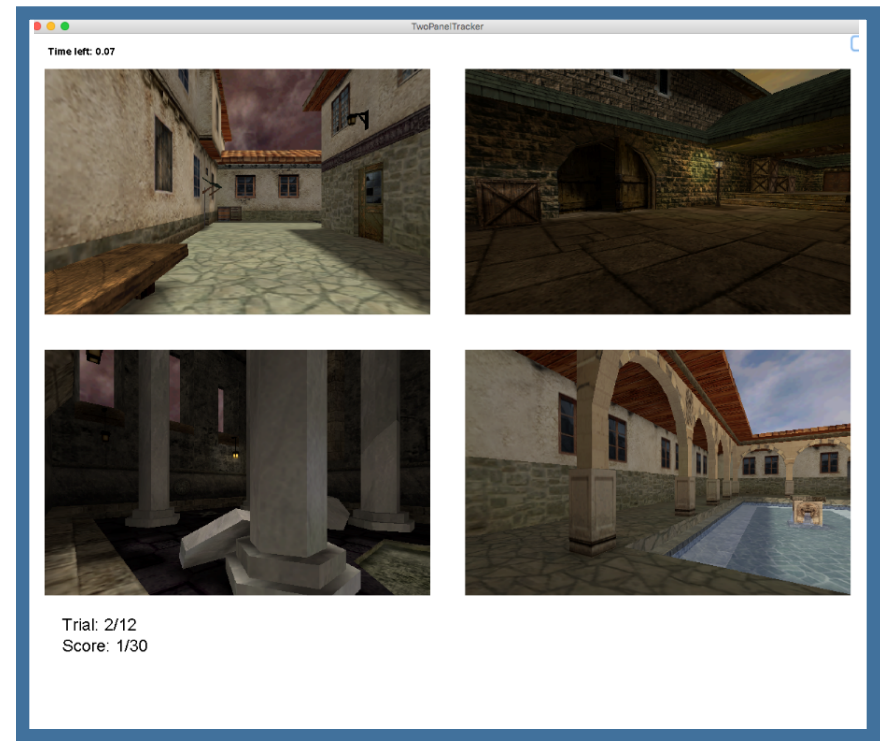
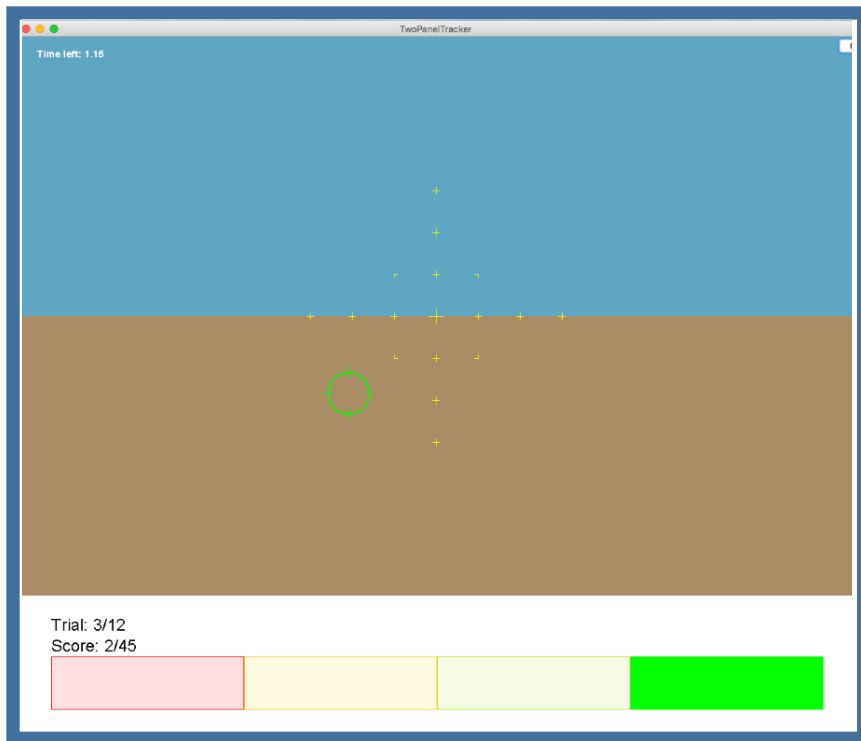


# Study 3 | Experiment setting

Evaluating Effects of User Experience and System Transparency on Trust in Automation

Thursday 1:30 pm

TOGGLE



Tracking Task



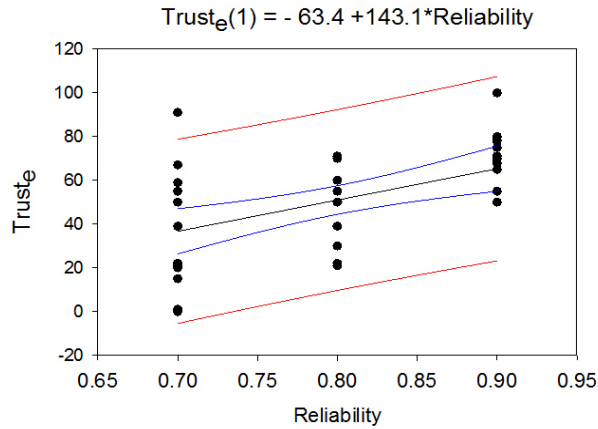
Detection Task



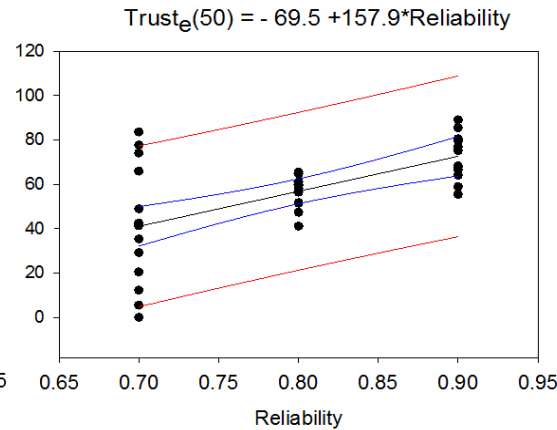
# Study 3 | Results

Binary Alarm

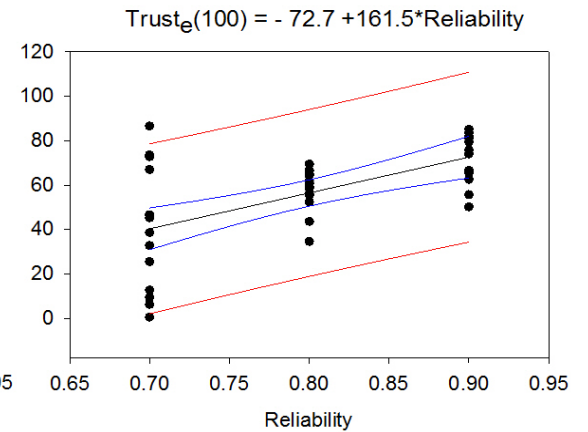
1<sup>st</sup> interaction



50<sup>th</sup> interaction

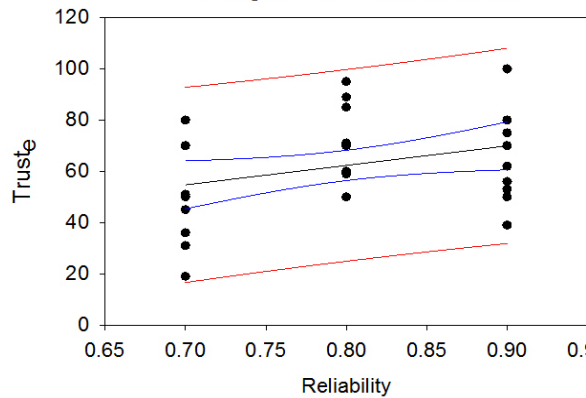


100<sup>th</sup> interaction

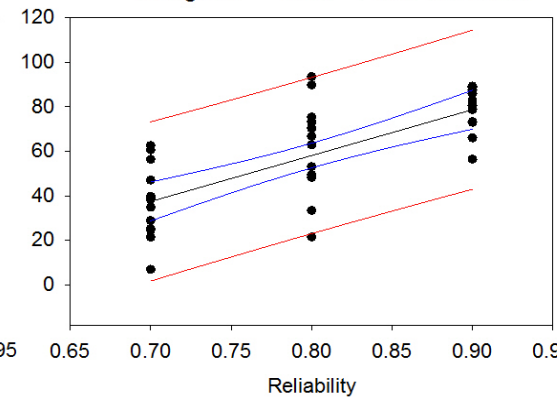


Likeli-  
Hood  
Alarm

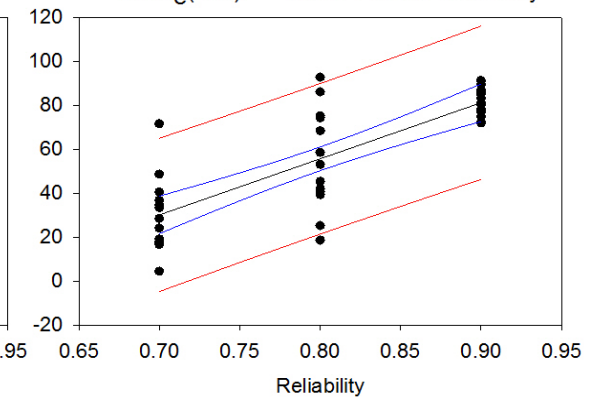
Trust<sub>e</sub>(1) = 1.4 + 76.2\*Reliability



Trust<sub>e</sub>(50) = -106.7 + 206.0\*Reliability



Trust<sub>e</sub>(100) = -148.6 + 255.4\*Reliability



- Over repeated interactions, users' trust becomes more appropriately calibrated
- With likelihood alarm, the calibration process is faster

# Takeaways

- Appropriate **calibration** is the key
- Focus on the characteristics of the automation and of the **human**
- To improve human-automation team performance
  - To minimize trust-reliability miscalibration
  - To minimize confidence-ability miscalibration



# Acknowledgement

