

4 July 2023 - Security for all in an AI enabled society

Cyber Hygiene in AI-enabled domestic life: A smart heating case study experiment

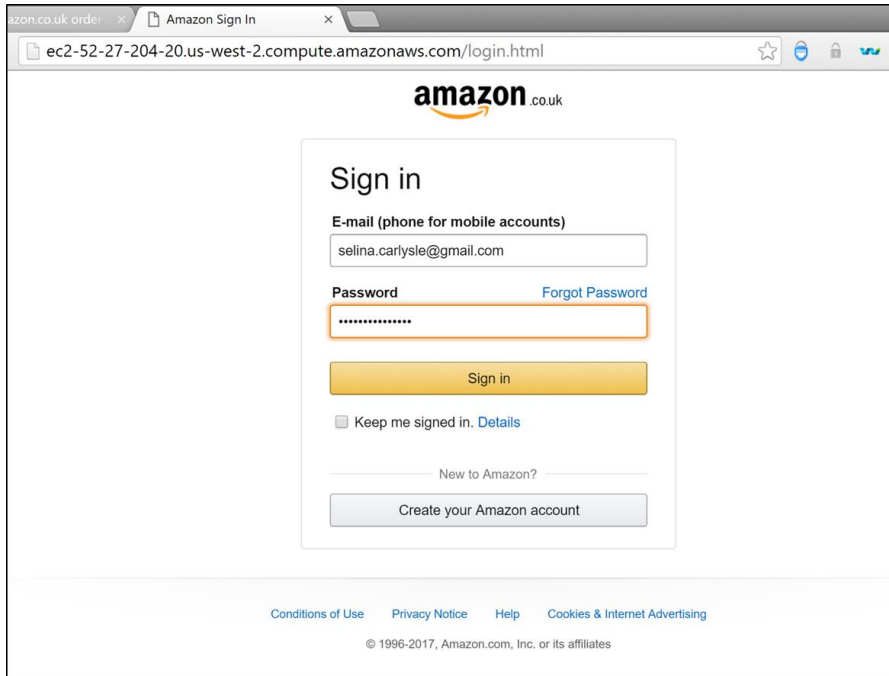
Professor George Loukas

University of Greenwich



From Human-as-a-Problem to Human-as-a-solution in cyber security

– “old-school (social media, web, email cyber threats)”



Commercial cyber security products Vs. humans

AntiVirus								Browsers								Platform			
Attack	A1	A2	A3	A4	A5	A6	A7	Attack	B1	B2	B3	B4	B5	B6	B7	Attack	P1	P2	P3
1.1	X	X	X	X	X	X	X	1.1	X	X	X	X	X	X	X	1.1	-	-	X
1.2	✓	X	X	X	X	X	X	1.2	✓	X	X	X	X	X	X	1.2	-	X	✓
2.1	X	X	X	X	X	X	X	2.1	X	X	X	X	X	X	X	2.1	X	-	X
2.2	X	X	X	X	X	X	X	2.2	X	X	X	X	X	X	X	2.2	X	-	X
3.1	X	X	X	X	X	X	X	3.1	X	X	X	X	X	X	X	3.1	-	-	X
3.2	X	X	X	X	X	X	X	3.2	X	X	X	X	X	X	X	3.2	-	-	X

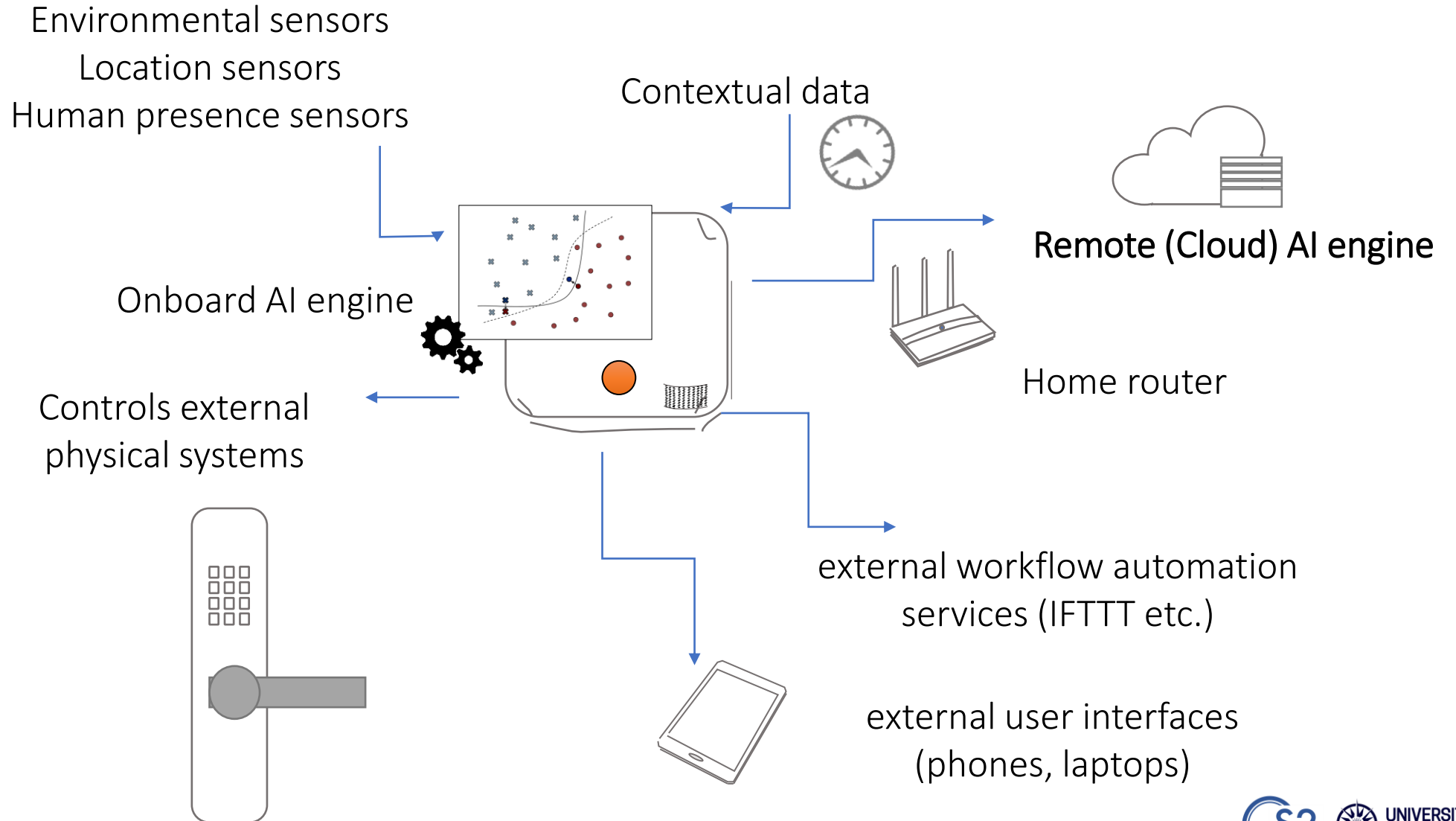
AntiVirus								Browsers								Platform		
Attack	A1	A2	A3	A4	A5	A6	A7	Attack	B1	B2	B3	B4	B5	B6	B7	Attack	P1	P3
2.1A	X	X	X	X	X	X	X	2.1A	X	X	X	X	X	X	X	2.1A	X	-
4.1	X	X	X	X	X	X	X	4.1	✓	X	X	X	X	X	X	4.1	X	-
4.2	X	X	X	X	X	X	X	4.2	X	X	X	X	X	X	X	4.2	X	-
5.1	X	X	X	X	X	X	X	5.1	-	-	-	-	-	-	-	5.1	-	X
5.2	✓	X	X	X	X	X	X	5.2	-	-	-	-	-	-	-	5.2	-	X

A	H1	H2	H3	H4	H5	H6	H7	H8	H9	H10	H11	H12	H13	H14	H15	H16	H17	H18	H19	H20	H21	H22	H23	H24	H25	H26
1.1	.49	.72	.54	.55	-	.59	X	.50	-	-	.70	.59	.48	-	.53	-	-	.33	.46	-	-	.46	.54	.70	.79	.66
1.2	-	.66	.53	-	-	-	✓	.26	-	-	-	-	.34	-	.18	-	-	.44	-	-	-	-	.17	-	.71	-
2.1A	-	-	-	-	-	.60	✓	-	-	-	.58	-	-	-	.48	-	-	-	-	.49	.30	-	-	-	-	-
2.1B	.51	-	-	-	-	-	✓	.51	-	-	.58	.43	-	-	.48	-	-	.32	.49	-	-	-	-	-	-	-
2.2	.51	-	-	-	-	-	-	-	-	-	-	-	-	.48	-	-	-	.32	-	-	-	-	-	-	-	-
3.1	-	-	.54	.55	-	.60	X	.63	-	.51	.70	.59	-	.56	.56	.46	.49	.33	.46	-	-	-	.54	.70	-	.66
3.2	-	-	-	-	-	-	.62	-	-	.49	-	-	-	.48	-	-	.26	.44	.46	-	-	-	-	.49	-	.67
4.1	-	-	-	.63	-	.56	-	-	-	.53	.72	.43	-	-	-	-	-	-	-	.49	-	-	-	-	-	-
4.2	-	-	-	-	-	-	-	-	-	.56	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
5.1	-	.67	-	.61	X	.63	X	.39	.35	-	.65	.52	-	-	.49	-	-	.45	.47	-	.30	.37	.66	.44	-	.68
5.2	-	-	-	-	X	-	✓	-	.35	-	-	-	-	-	-	-	-	-	.47	-	.57	-	.68	.72	-	.64

2 in 3 get it right

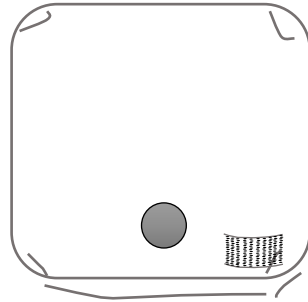
Heartfield, R. and Loukas, G. (2018) Detecting semantic social engineering attacks with the weakest link. Computers & Security, Elsevier.

Is the concept still applicable in AI-enabled society?



Is the concept still applicable in AI-enabled society?

“It’s playing up again”



Is the concept still applicable in AI-enabled society?

There are many reasons to want it to be.

- People **feeling more in control** of their own devices and systems
- **Extra layer of defence** especially for new and unknown threats to AI
- People **will not trust AI blindly**

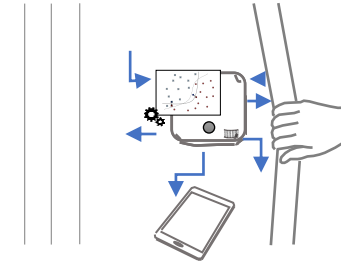
Is the concept still applicable in AI-enabled society?

In CHAI, we explore this from four angles

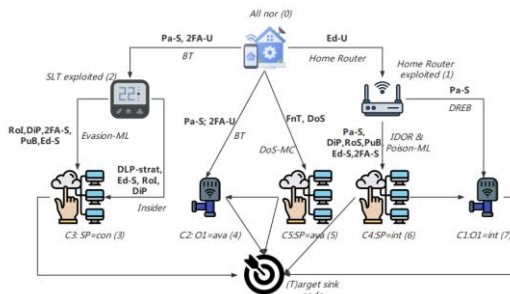
Cyber hygiene for AI users (+ Training)

Perception (What to notice)
Detection (How to diagnose)
Response (How to act)

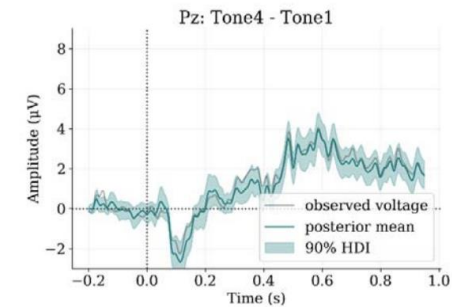
AI transparency for non-experts



Optimisation of AI measures for the service provider



Experimental study of the neuroscience dimensions



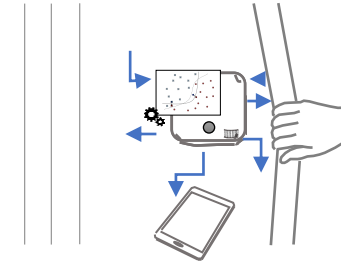
Is the concept still applicable in AI-enabled society?

The first two are linked directly

Cyber hygiene for AI users (+ Training)

Perception (What to notice)
Detection (How to diagnose)
Response (How to act)

AI transparency for non-experts



Cyber hygiene

Cyber hygiene clearly leads to online cyber safety, encompassing what we think, how we think, and what we do online

Vishwanath, A., Neo, L.S., Goh, P., Lee, S., Khader, M., Ong, G. and Chin, J., 2020. Cyber hygiene: The concept, its measure, and its initial tests. *Decision Support Systems*, 128.

In cybersecurity, humans are considered a problem to be controlled. This robs organisations of their human agents' ability to make a contribution to cyber security.

Zimmermann, V. and Renaud, K., 2019. Moving from a 'human-as-problem' to a 'human-as-solution' cybersecurity mindset. *International Journal of Human-Computer Studies*, 131, pp.169-187.

AI transparency for non-experts

Gaining even limited knowledge of the inner workings of AI approaches is important for non-experts to develop trust.

Ayobi, A., Stawarz, K., Katz, D., Marshall, P., Yamagata, T., Santos-Rodríguez, R., Flach, P. and O'Kane, A.A., 2021. Machine learning explanations as boundary objects: how AI researchers explain and non-experts perceive machine learning.

AI risk communication suffers from disconnect between practitioners' and researchers' mental models of threats.

Bieringer, L., Grosse, K., Backes, M., Biggio, B. and Krombholz, K., 2022. Industrial practitioners' mental models of adversarial machine learning. In *SOUPS 2022*.

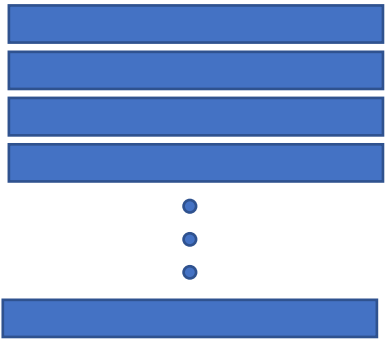
Explainable AI can help secure human-interactive robots

Roque, A. and Damodaran, S.K., 2022. Explainable AI for Security of Human-Interactive Robots. *International Journal of Human-Computer Interaction*, 38(18-20), pp.1789-1807.

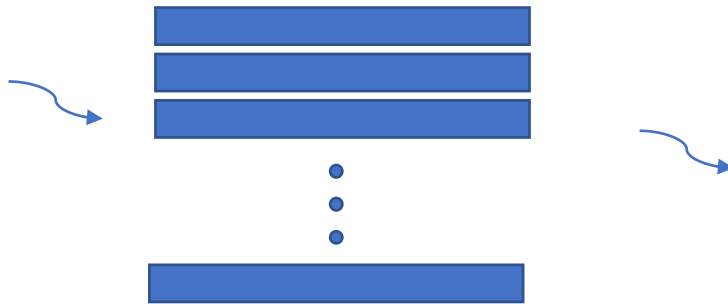
Existing Cyber Hygiene and AI-enabled smart homes

Analysed applicability of 105 IoT security recommendations provided by NCSC, ENISA, etc. in terms of **applicability to prevention, perception, detection and response for AI-IoT, population** (expertise required), **intervention** (time required), and **scientific support**.

All recommendations
105



Applicable for common AI-IoT
with no heavy expertise or time
75



Substantial scientific evidence
only 5
(all preventive and only indirectly useful)

- PR-1: Turn off service when not needed
- PR-2: Change all default passwords or PINs
- PR-3: Use different passwords for different systems
- PR-4: Do not use passwords that are based on personal information that can be easily accessed or guessed
- PR-5: Do not use passwords that can be found in any dictionary of any language

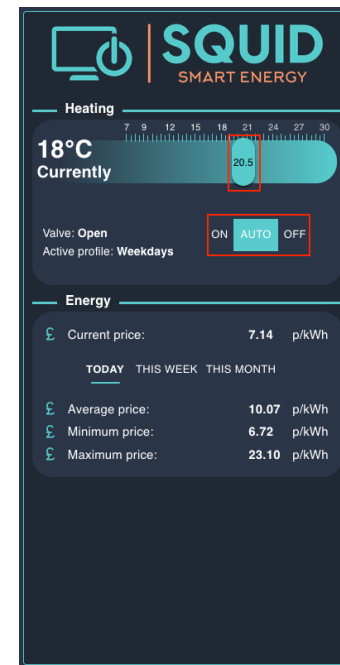
A smart heating case study (AI-controlled radiator valve)

Aiming for simplest possible scenario to explain to a non-expert:

Schedule learning for **radiator valve setpoint**
Vs. price at different times of the day.

Bayesian linear regression

Separate model for each of five profiles
(mornings, weekdays/weekends, evenings, nights)



<https://github.com/chai-project>

AI Transparency measures in Squid application (1/2)

Aim: To provide more context on AI decisions than is usually available to the users, so as to help them notice AI misbehaviour

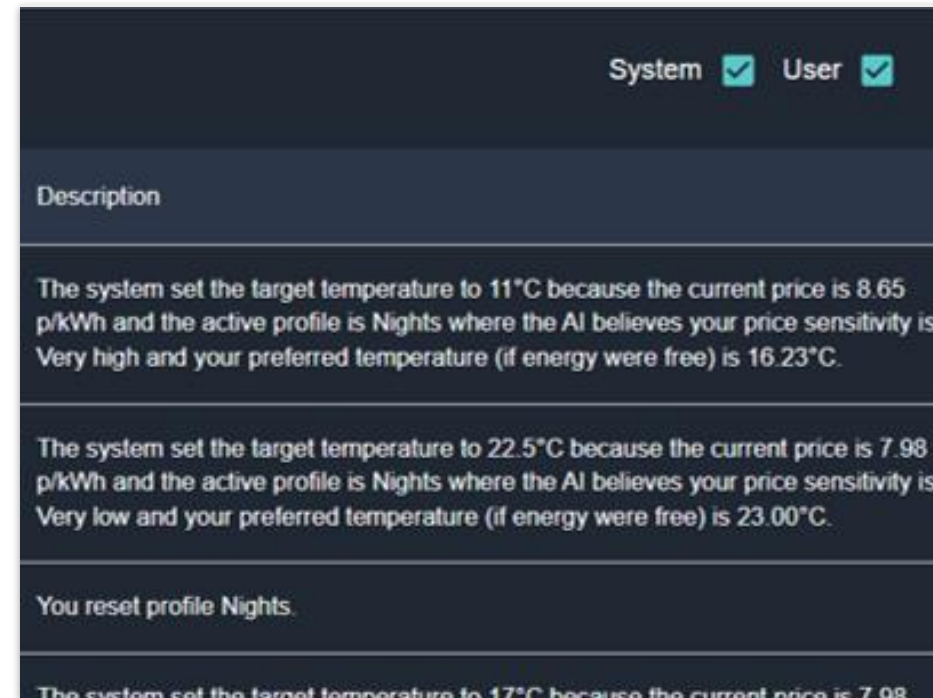
Preferred temp at price 0 and price sensitivity



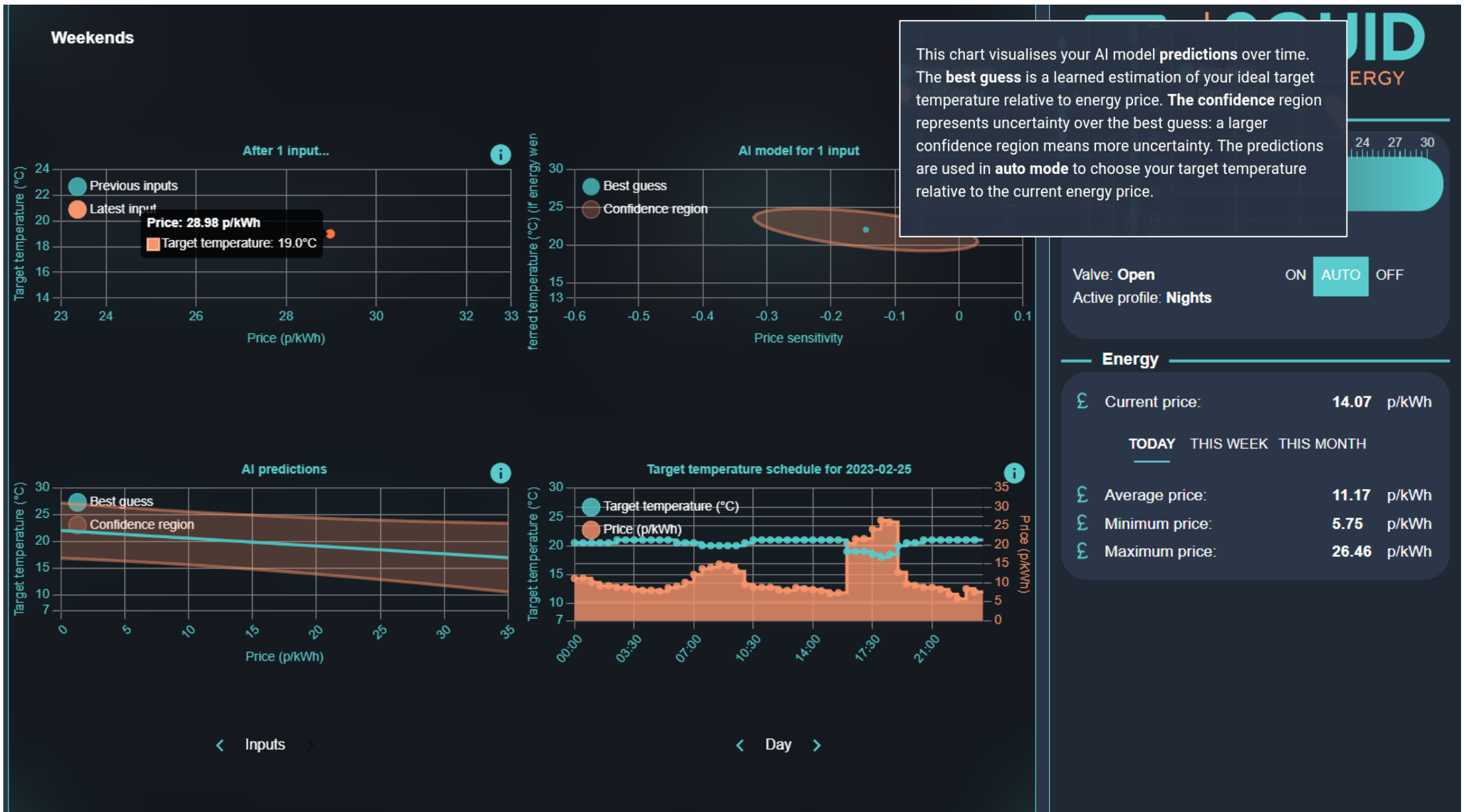
Log of notifications

Example:

*“The system set the target temperature to **11 C** because the current price is **8.65 p/kWh** and the active profile is **Nights** where the AI believes your price sensitivity is **Very high** and your preferred temperature (if energy were free) is **16.23 C.**”*



AI Transparency measures in Squid application (2/2)



This chart visualises your AI model **predictions** over time. The **best guess** is a learned estimation of your ideal target temperature relative to energy price. The **confidence** region represents uncertainty over the best guess: a larger confidence region means more uncertainty. The predictions are used in **auto mode** to choose your target temperature relative to the current energy price.

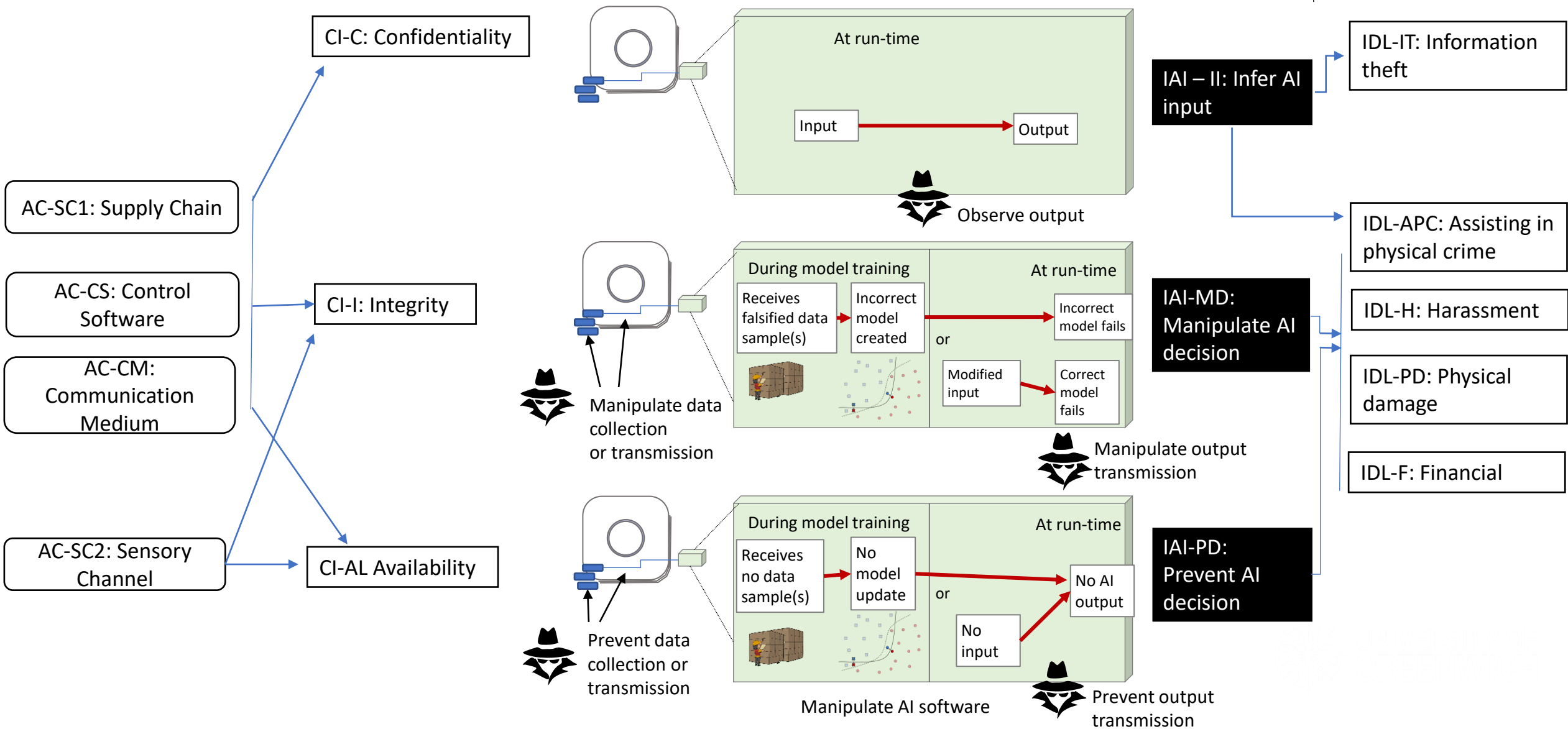
Taxonomy of attacks on AI-enabled smart homes

AV: Attack vector

CI: Cyber impact

IAI: Impact on AI operation

IDL: Impact on domestic life



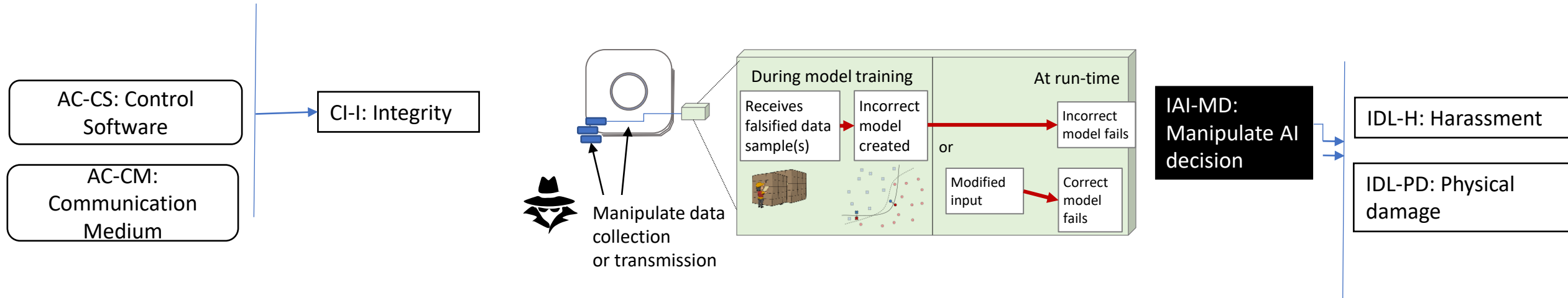
Attacks in CHAI experiment

AV: Attack vector

CI: Cyber impact

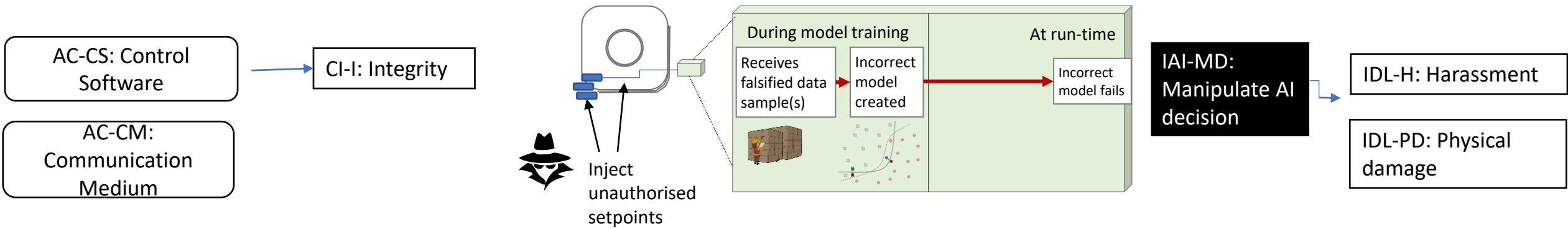
IAI: Impact on AI operation

IDL: Impact on domestic life



Focused on integrity attacks because they leave physical and digital traces observable by non-experts

Attacks in CHAI experiment (1/3)



Simple poisoning to modify price sensitivity or preferred temp at price 0.

Method: Unauthorised setpoint injections

From: 1/4/2022 To: 25/4/2022 System User

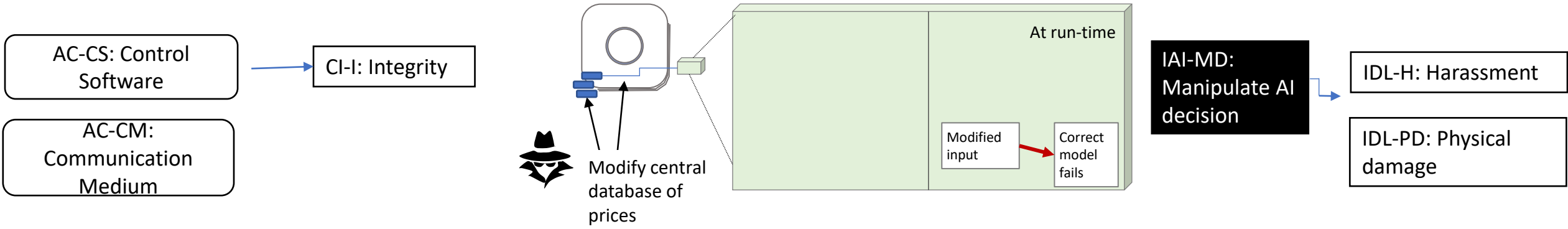
Date	Time	Category	Description
5/4/2022	08:20	User	You set the target temperature to 7 C (override mode is now active).
5/4/2022	08:21	User	You set the target temperature to 21 C (override mode is now active).
5/4/2022	08:22	User	You set the target temperature to 7 C (override mode is now active).
5/4/2022	08:23	User	You set the target temperature to 21 C (override mode is now active).
5/4/2022	08:24	User	You set the target temperature to 7 C (override mode is now active).
5/4/2022	08:25	User	You set the target temperature to 21 C (override mode is now active).
5/4/2022	08:26	User	You set the target temperature to 7 C (override mode is now active).
5/4/2022	08:27	User	You set the target temperature to 21 C (override mode is now active).
5/4/2022	08:28	User	You set the target temperature to 7 C (override mode is now active).
5/4/2022	08:29	User	You set the target temperature to 21 C (override mode is now active).
5/4/2022	08:30	User	You set the target temperature to 7 C (override mode is now active).

Rows per page: 25 26-50 of 378 < >

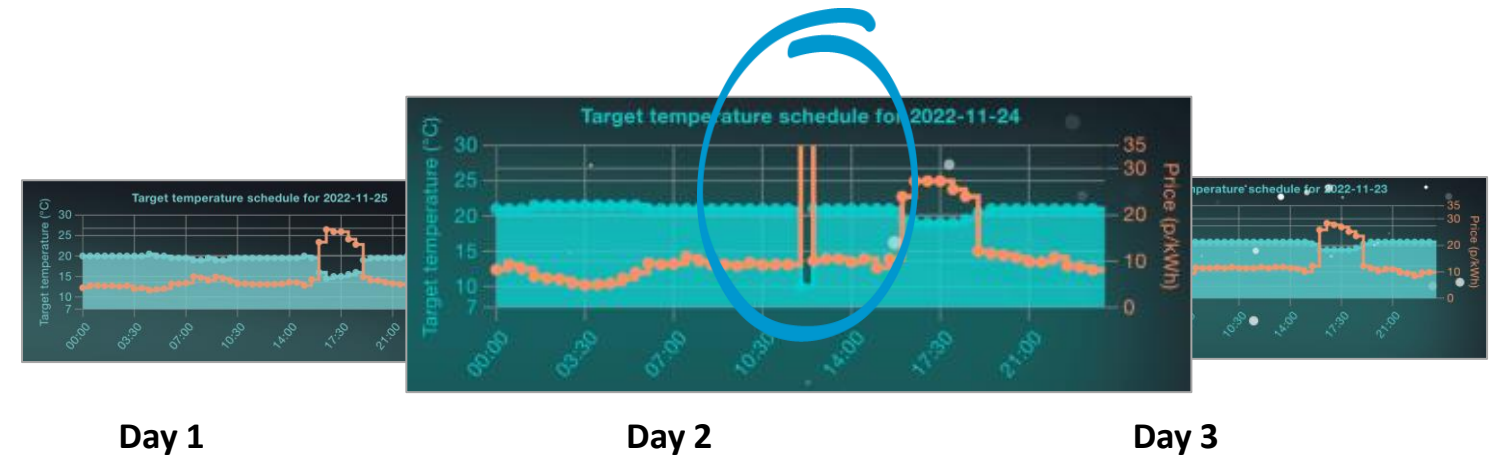
HOME SCHEDULE PROFILES NOTIFICATIONS



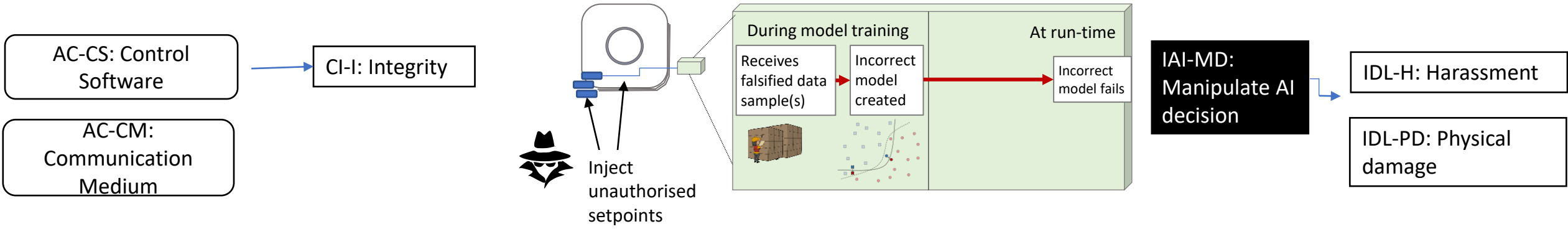
Attacks in CHAI experiment (2/3)



Price manipulation to affect AI setpoints at scale
Method: Central database modified for a period of time



Attacks in CHAI experiment (1/3)



Stealthy poisoning to modify price sensitivity or preferred temp at price 0.
Method: Unauthorised setpoint injections and deletion of log entries

Date	Time	Description
05/12/2022	10:59	The system set the target temperature to 11°C because the current price is 8.65 p/kWh and the active profile is Nights where the AI believes your price sensitivity is Very high and your preferred temperature (if energy were free) is 16.23°C.
05/12/2022	10:59	The system set the target temperature to 22.5°C because the current price is 7.98 p/kWh and the active profile is Nights where the AI believes your price sensitivity is Very low and your preferred temperature (if energy were free) is 23.00°C.
05/12/2022	10:59	You reset profile Nights.
05/12/2022	10:59	The system set the target temperature to 17°C because the current price is 7.98 p/kWh and the active profile is Nights where the AI believes your price sensitivity is High and your preferred temperature (if energy were free) is 19.76°C.
05/12/2022	10:59	The system set the target temperature to 16.5°C because the current price is 7.98 p/kWh and the active profile is Nights where the AI believes your price sensitivity is High and your preferred temperature (if energy were free) is 19.76°C.



Cyber hygiene for perception, detection and response

A simple diagnostic graph method

Truth table of observable indicators -> Binary decision diagram for non-experts

	Diagnostics (persistent facts)				Perception (temporary observations that trigger suspicion)		
	D1	D2	D3	D4	P1	P2	P3
Simple poisoning	0	0	0	1	1	1	1
Input manipulation	1	1	0	0	0	0	0
Stealthy poisoning	1	0	0	1	1	1	1
<i>Normal rapid setpoint changes</i>	1	0	1	1	1	1	1

D1: all user entries are recognised

D2: dramatic price change

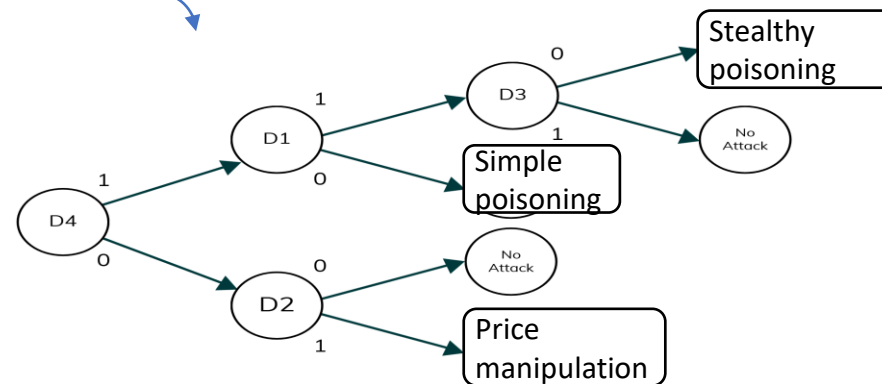
D3: several recent user entries

D4: significant preferred temperature at price 0 change

P1: frequent buzz

P2: room too warm/cold

P3: unusual price sensitivity



Cyber hygiene for perception, detection and response

Diagnostic graph UI and indicative content



HELP | **SQUID**
SMART ENERGY

This tool helps you diagnose possible cyberattacks against the Squid app.

You will first be asked about what made you suspicious of a cyberattack. Then we will help you figure out what you can do about it.

- I have heard the buzzing sounds from the smart valve more than **three times** per hour
- The room feels **too cold** for this time of day
- On the gauge the needle **moved unusually**: if the needle is pointing to the **red sections** or if it has moved **two positions or more**, the AI could be acting unexpectedly.
- Something else

Go to NOTIFICATIONS

In the description column, check the estimated "preferred temperature (if energy were free)". Has it changed by more than **2°C** in the last 2 days?

! A security attack would aim for a large change in temperature.

Time	Temp	Category	Description
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?

Yes

No

Go to NOTIFICATIONS

Check the most recent log entries under the 'user' category. Look at the date and time columns. Do you remember making these temperature changes?

! If you did not make these changes, someone could be adding data without your permission.

Time	Temp	Category	Description
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?
2023-01-01 10:00	18.0	user	The system set the target temperature to 17°C because the user profile is 1.80 and the valve profile is higher when the AI balances your preferences. Has it changed by more than 2°C in the last 2 days?

Yes

No

You may have detected an AI Poisoning Attack

! It seems that a number of temperature changes have been entered into the system over a short period of time

These entries have become part of your current profile. Squid is heating the radiator according to this new data

WHAT TO DO NEXT

- ✓ Look at the problematic log description and **identify which profile has been breached**. This needs to be **reset**
- ✓ Go to *Profiles*. Select the name of the **breached profile** from the top left menu and **click on the Reset this Profile** button on the top right-hand side

Cyber hygiene for perception, detection and response

Experimental setup

19 participants from 10 households

Schedule:

Training on Squid

1 month of familiarisation

Training on one attack

1 month of attacks (2 x 3 attacks each)

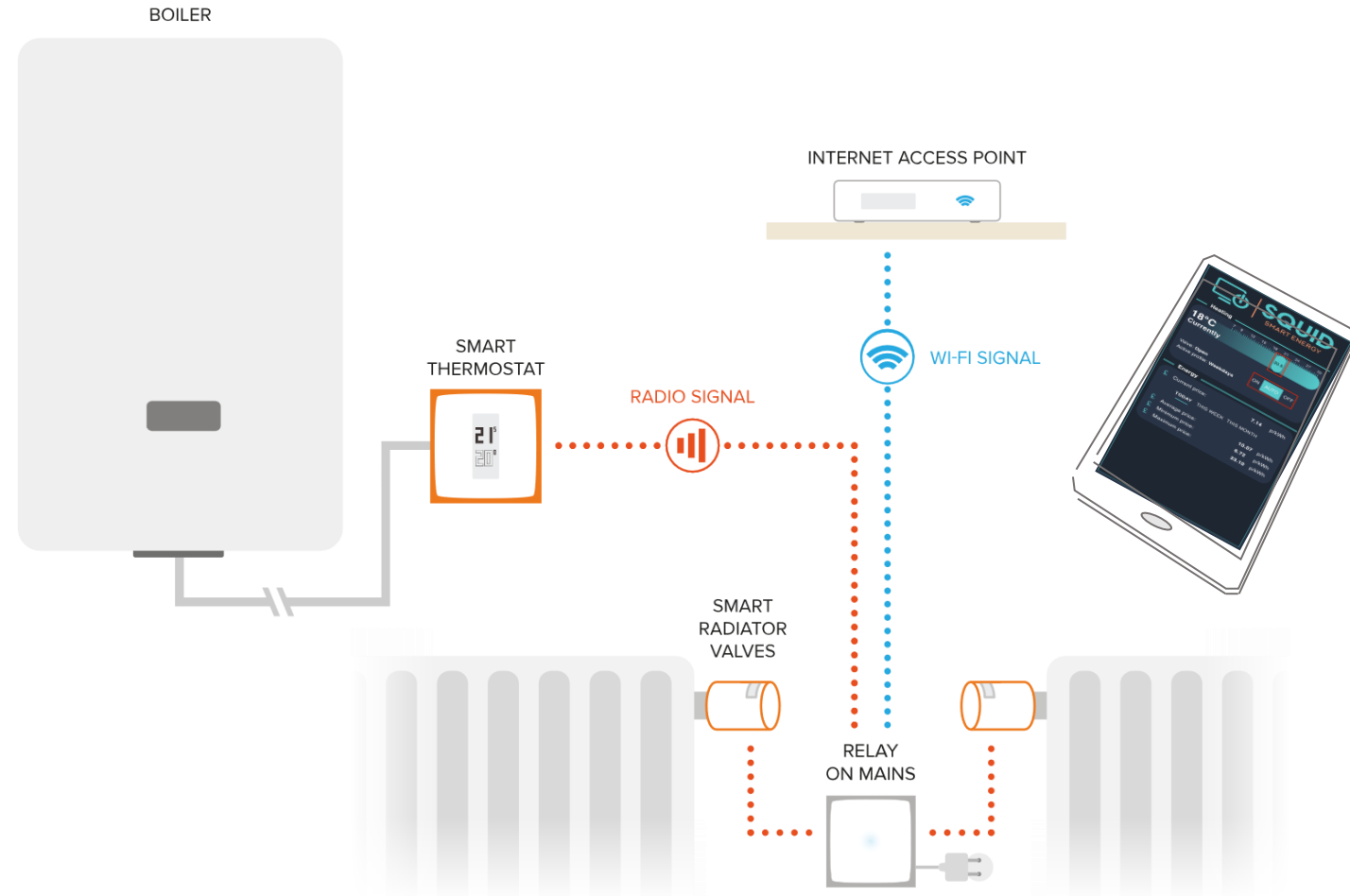
Data collection:

Reporting diary for each incident suspected

Weekly interviews with each household

Squid usage data

Diagnostic tool usage data



Cyber hygiene for perception, detection and response Training

1 hour session including slides and use of the diagnostic tool on the Simple poisoning attack.
No training on the AI transparency features or other attacks

Perception

What signs did you identify?

APP BEHAVIOUR

The target temperature on the app has **changed** a lot in the past 30 mins

ENVIRONMENT

The room feels a lot **warmer** than I would expect

DEVICE FUNCTION

After increasing the target temp it took a **long time** for the valve to turn on

SOUNDS

I heard the radiator valve **clicking** on/off more often than normal

APP BEHAVIOUR

The target temperature on the app has **not changed** recently

APP BEHAVIOUR

The **AI logs** do not appear to be correct

SOUNDS

There is **no sound** from the radiator valve when the heating turns on

Detection and response



Troubleshooting your Squid app. Do you observe any of these?

Any first observations of yours can help narrow down what may have happened. Tick as many as you observed

	Responses		
	Yes	No	Do not know
I have heard the radiator valve activating very frequently or very infrequently.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
The target temperature shown on the app does not seem to have changed recently.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
The target temperature on the app has changed many times in the last half-hour.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
I noticed a long delay between the heating being requested via a change in target temperature and the radiator valve actually activating.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
The AI logs do not appear to be correct.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
Cannot hear any sound from the radiator valve when the heating is required.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
The radiator/room feels too warm or too cold to what you would normally expect.	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>

Cyber hygiene for perception, detection and response

Experimental results 1 – Human as an AI threat sensor

Perception

Groundtruth		
	Positive	Negative
Positive	47 (78%)	8 (4%)
Negative	13 (22 %)	212* (96%)

f1 score = 0.86

Detection

Groundtruth		
	Positive	Negative
Positive	31 (60%)	2 (25%)
Negative	13 (40%)	6 (75%)

f1 score = 0.73

Response

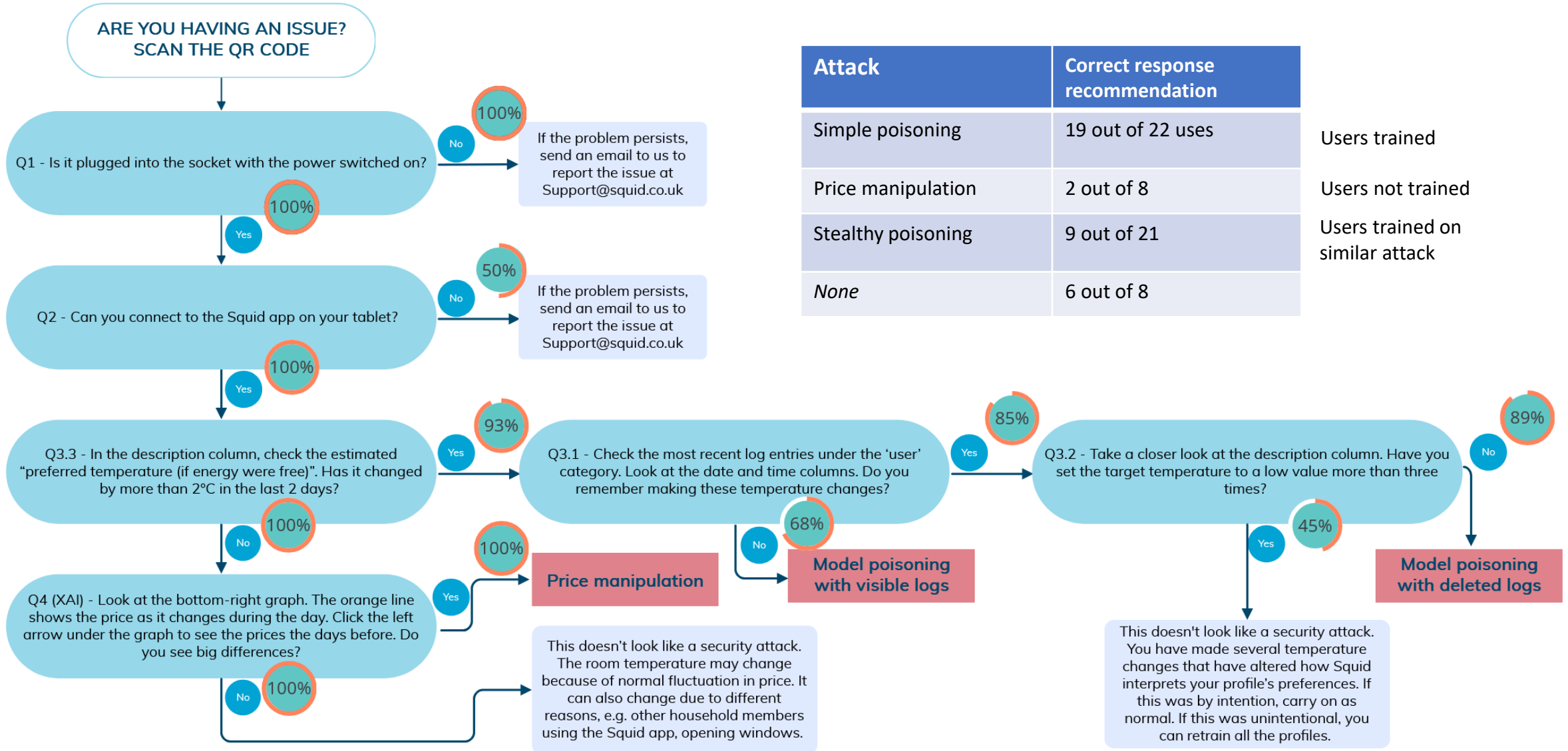
Groundtruth		
	Attack	Non-attack
Correct Action	39 (78%)	6 (86%)
Incorrect Action	11 (22%)	1 (14%)

Also, no obvious correlation between level of engagement (number of interactions) and performance as human sensors

* True negative is a day where non-attack occurred and no attack was perceived

Cyber hygiene for perception, detection and response

Experimental results 3 – diagnostic tool



Initial insights

Existing cyber hygiene for AI-IoT is only preventive and not AI-specific.

We developed a domain-specific example set of perception, detection and response guidance, supported by partial training.

Perception accuracy was high (f1 score 0.86)

Detection worked well only for the simple attack we had already trained the participants. Much less so for the rest.

No correlation between number of interactions and human sensor performance

and next step

Test these initial insights quantitatively in lab-based experiment

Compare different AI explainability approaches

Assess practicality on different AI application areas

Thank you