# A Review of AI Agent Reasoning with Values
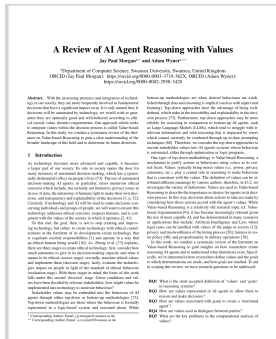
Jay Paul Morgan & Adam Wyner

**Computational** Foundry
Ffowndri **Gyfrifiadol**

Swansea University
Prifysgol **Abertawe**

25th October 2025
Value Engineering in AI (VALE) 2025

What we have done:

- Performed a systematic review of the literature on Value-based Reasoning.
- Articles for consideration were selected using a standardised (systematic) process.
- From the 57 articles in question, information was extracted to answer 5 research questions.

In this talk:

1. A brief description of the systematic framework.
2. 3 of the research questions and what has been found.
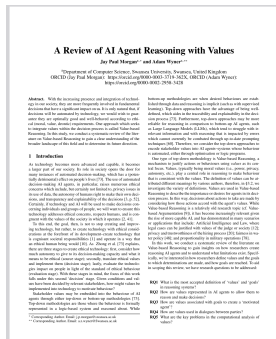3. Conclusions & future direction.

What we have done:

- Performed a systematic review of the literature on Value-based Reasoning.
- Articles for consideration were selected using a standardised (systematic) process.
- From the 57 articles in question, information was extracted to answer 5 research questions.

In this talk:

1. A brief description of the systematic framework.
2. 3 of the research questions and what has been found.
3. Conclusions & future direction.

Table: PICO keywords and synonyms used to search digital libraries.

|            | **Keywords**                                                          | **Synonyms**                                                          |
|------------|-----------------------------------------------------------------------|----------------------------------------------------------------------|
| Population | AI                                                                    | AI Agents, Artificial Intelligence, Machine Learning, Multi-Agent Systems |
| Intervention | Value-based Reasoning                                               | Value-based argumentation                                            |
| Comparison | Logic, Argumentation, Reasoning, Philosophy, Value, Computational Model | Norm                                                                |
| Outcome    | Representation, Dialogue, Behaviour                                    | Negotiation, Persuasion, Action                                      |

Figure: Overview of the article selection process.

(a) The number of articles selected from 2000-2025.

(b) The top-10 sources from which articles were selected.

# Research Questions

## RQ1

What is the most accepted definition of 'values' and 'goals' in reasoning systems?

## RQ2

How are values represented in AI agents to allow them to reason and make decisions?

## RQ3

How are values associated with goals to create a 'motivated agent'?

## RQ4

How are values used in dialogues between parties?

## RQ5

What are the key problems in the computational analysis of values?

**RQ1**: What is the most accepted definition of 'values' and 'goals' in reasoning systems?

How do different authors define 'values' and 'goals' and do these have an effect on the formalisations? And is there a consensus?

Values:

- Values are "Abstract principles that guide behaviour"
- Many look to the Schwartz's Theory of Basic Human Values (STBHV).
- When Schwartz is used, sometimes only a some of values are used–which limit the usefulness of STBHV.

Goals:

- Tend to be less defined.
- "Goals reflect the state of affairs the agent wishes to bring about".
- Some types of Goals:
    - Achievement (Make false → true)
    - Remedy (Make true → false)
    - Maintenance (Keep true, true)
    - Avoidance (Keep false, false)

## RQ1

What is the most accepted definition of 'values' and 'goals' in reasoning systems?

## RQ2

How are values represented in AI agents to allow them to reason and make decisions?

## RQ3

How are values associated with goals to create a 'motivated agent'?

## RQ4

How are values used in dialogues between parties?

## RQ5

What are the key problems in the computational analysis of values?

**RQ2**: How are values represented in AI agents to allow them to reason and make decisions?

Throughout the articles, we see two methodologies for creating such reasoning systems:

(1) through value-alignment, meaning values are implicitly encoded into the system through its output behaviour;

(2) an explicit representation of values where determinations of inner behaviour and action are reasoned through states/functions that represent the values of interest.

**Implicit**
- No explicit representation.
- System performs instrumental actions (actions that are instrumental to bring out goals which align with values).

**Explicit**
- State/Object
  $V = \{v_1, v_2, ..., v_n\}$
- Function
  $f(\text{action})_v \rightarrow \{+, -, =\}$
- Numerical
  $v_1 = 0.25, v_2 = 0.5, ...$

## RQ1

What is the most accepted definition of 'values' and 'goals' in reasoning systems?

## RQ2

How are values represented in AI agents to allow them to reason and make decisions?

## RQ3

How are values associated with goals to create a 'motivated agent'?

## RQ4

How are values used in dialogues between parties?

## RQ5

What are the key problems in the computational analysis of values?

Lastly, we evaluate the key issues with the formalisations that might point to the future direction of the field.

- **Deeper and more realistic scenarios** Work should be done to create datasets or scenarios with which to compare methodologies.
- Connection between **Values & Goals**. While 'value' has had more attention in definition, the definition of 'goal' remains somewhat implicit (with only a few articles giving some definition and types). Furthermore, the connection between short-term and longer-term goals and how these are strategised with values can be explored in more depth (**Planning**).
- Concept of **Preferences**. Preferences have been limited to ordering relations, but this concept could be taken a lot further.

- Conducted a literature review that has identified main themes such as a general direction on the meaning of Values.

- Some concepts can be taken further: Connection between Values & Goals; Preferences.

- This presentation has only touched upon the main points, but there is more. Come and ask some questions at the poster session!

- Conducted a literature review that has identified main themes such as a general direction on the meaning of Values.

- Some concepts can be taken further: Connection between Values & Goals; Preferences.

- This presentation has only touched upon the main points, but there is more. Come and ask some questions at the poster session!