

# Truth-Tracking with Non-Expert Information Sources

---

Joe Singleton and Richard Booth

singletonj1@cardiff.ac.uk

- **Problem:** what can we learn from non-expert information sources?
- We aim to learn both:
  - the **true facts** of the world
  - the **true level of expertise** of the sources
- We adapt the learning framework from recent work combining formal learning theory, belief revision and epistemic logic
- Main results:
  - description of **what can be learned**
  - characterisation of **truth-tracking learning methods**

- **Problem:** what can we learn from non-expert information sources?
- We aim to learn both:
  - the **true facts** of the world
  - the **true level of expertise** of the sources
- We adapt the learning framework from recent work combining formal learning theory, belief revision and epistemic logic
- Main results:
  - description of **what can be learned**
  - characterisation of **truth-tracking learning methods**
- **Warning:** Still preliminary work! Strong assumptions on the input the learning method receives ⚠

## Motivating example

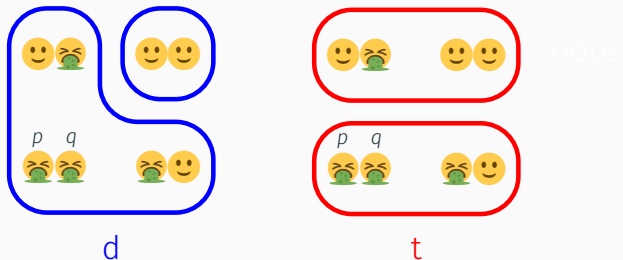
---

## Motivating example

- Patient  $a$  is checked for conditions  $p$  and  $q$
- Doctor  $d$  has expertise to determine whether whether  $a$  has at least one condition, but needs a blood test to tell which one(s)
- Tests are only available for  $p$ : tech  $t$  has expertise on  $p$  but not  $q$

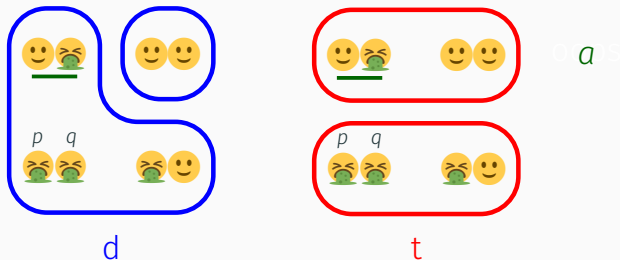
## Motivating example

- Patient  $a$  is checked for conditions  $p$  and  $q$
- Doctor  $d$  has expertise to determine whether whether  $a$  has at least one condition, but needs a blood test to tell which one(s)
- Tests are only available for  $p$ : tech  $t$  has expertise on  $p$  but not  $q$
- We model expertise with **partitions** of states: sources cannot distinguish states in the same cell



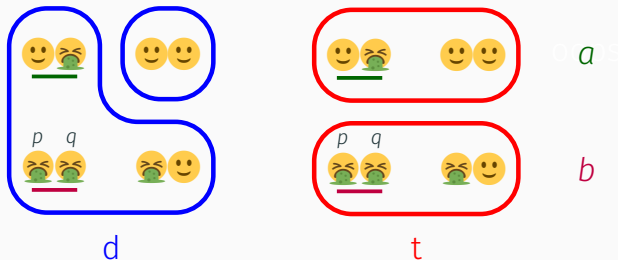
## Motivating example

- Patient  $a$  is checked for conditions  $p$  and  $q$
- Doctor  $d$  has expertise to determine whether whether  $a$  has at least one condition, but needs a blood test to tell which one(s)
- Tests are only available for  $p$ : tech  $t$  has expertise on  $p$  but not  $q$
- We model expertise with **partitions** of states: sources cannot distinguish states in the same cell



## Motivating example

- Patient  $a$  is checked for conditions  $p$  and  $q$
- Doctor  $d$  has expertise to determine whether whether  $a$  has at least one condition, but needs a blood test to tell which one(s)
- Tests are only available for  $p$ : tech  $t$  has expertise on  $p$  but not  $q$
- We model expertise with **partitions** of states: sources cannot distinguish states in the same cell





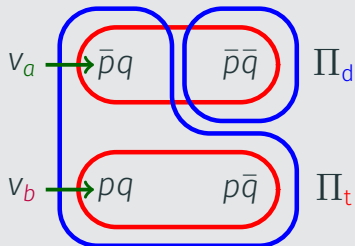
## Logical framework for expertise

---

## Basic framework

- $\mathcal{P}$ : finite set of propositional variables (e.g.  $p, q, \dots$ )
- $\mathcal{S}$ : finite set of sources (e.g.  $d, t, \dots$ )
- $\mathcal{C}$ : finite set of cases (e.g.  $a, b, \dots$ )
- Valuation:  $v : \mathcal{P} \rightarrow \{0, 1\}$
- A **world** is a pair  $W = (\{\Pi_i\}_{i \in \mathcal{S}}, \{v_c\}_{c \in \mathcal{C}})$ , where
  - Each  $\Pi_i$  is a partition of the set of valuations
  - Each  $v_c$  is a valuation

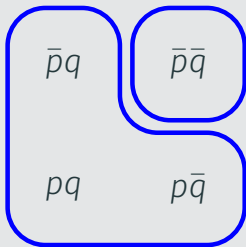
### Example



## Expertise and permissibility

- *i* **has expertise** on  $\varphi$  if *i* can always determine the correct value of  $\varphi$ :  
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$
- $\varphi$  states always distinguishable from  $\neg\varphi$  states

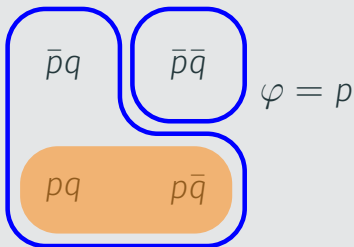
### Example



## Expertise and permissibility

- $i$  **has expertise** on  $\varphi$  if  $i$  can always determine the correct value of  $\varphi$ :  
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$
- $\varphi$  states always distinguishable from  $\neg\varphi$  states

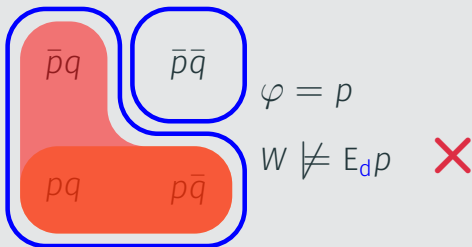
### Example



## Expertise and permissibility

- $i$  has expertise on  $\varphi$  if  $i$  can always determine the correct value of  $\varphi$ :  
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$
- $\varphi$  states always distinguishable from  $\neg\varphi$  states

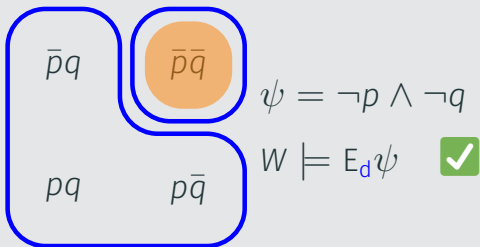
### Example



## Expertise and permissibility

- $i$  **has expertise** on  $\varphi$  if  $i$  can always determine the correct value of  $\varphi$ :  
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$
- $\varphi$  states always distinguishable from  $\neg\varphi$  states

### Example



## Expertise and permissibility

- *i* **has expertise** on  $\varphi$  if *i* can always determine the correct value of  $\varphi$ :

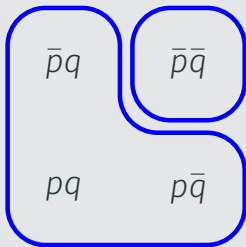
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$

- $\varphi$  states always distinguishable from  $\neg\varphi$  states
- $\varphi$  is **permissible** for *i* if  $\varphi$  is true up to lack of expertise of *i*

$$W, c \models P_i \varphi \iff \Pi_i[v_c] \cap \text{mods}(\varphi) \neq \emptyset$$

- true state indistinguishable from some  $\varphi$  state

### Example



## Expertise and permissibility

- *i* **has expertise** on  $\varphi$  if *i* can always determine the correct value of  $\varphi$ :

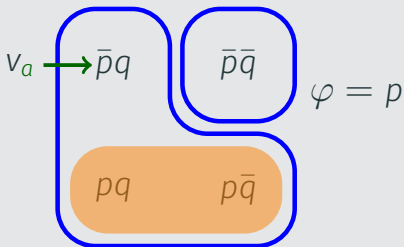
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$

- $\varphi$  states always distinguishable from  $\neg\varphi$  states
- $\varphi$  is **permissible** for *i* if  $\varphi$  is true up to lack of expertise of *i*

$$W, c \models P_i \varphi \iff \Pi_i[v_c] \cap \text{mods}(\varphi) \neq \emptyset$$

- true state indistinguishable from some  $\varphi$  state

### Example





## Expertise and permissibility

- $i$  **has expertise** on  $\varphi$  if  $i$  can always determine the correct value of  $\varphi$ :

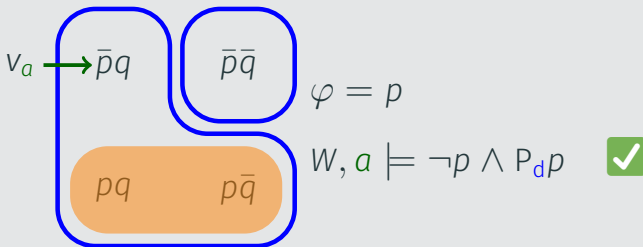
$$W \models E_i \varphi \iff (u \in \text{mods}(\varphi) \implies \Pi_i[u] \subseteq \text{mods}(\varphi))$$

- $\varphi$  states always distinguishable from  $\neg\varphi$  states
- $\varphi$  is **permissible** for  $i$  if  $\varphi$  is true up to lack of expertise of  $i$

$$W, c \models P_i \varphi \iff \Pi_i[v_c] \cap \text{mods}(\varphi) \neq \emptyset$$

- true state indistinguishable from some  $\varphi$  state

### Example

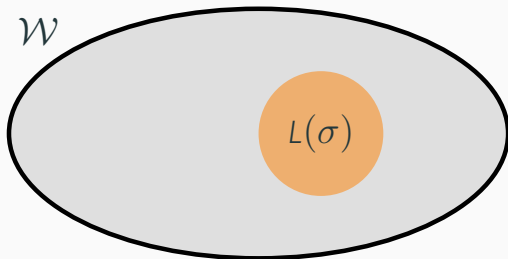


## Learning and truth-tracking

---

## Reports and methods

- We receive **reports** of the form  $\langle i, c, \varphi \rangle$ 
  - “source  $i$  reports  $\varphi$  in case  $c$ ”
- A **learning method**  $L$  maps a finite sequence  $\sigma$  to a **conjecture**  $L(\sigma) \subseteq \mathcal{W}$ , where  $\mathcal{W}$  is the set of all worlds



### Example

$$L(\sigma) = \{W \mid \forall \langle i, c, \varphi \rangle \in \sigma, W, c \models P_i \varphi\}$$

- We assume sources report all they consider possible
  - All reports are **permissible**: only false due to lack of expertise
  - All permissible reports **eventually appear**

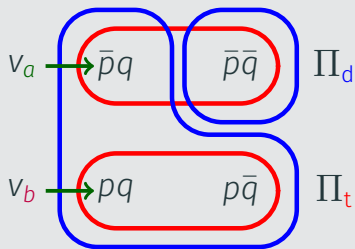
- We assume sources report all they consider possible
  - All reports are **permissible**: only false due to lack of expertise
  - All permissible reports **eventually appear**
- **Warning**: Strong assumptions! Sources are always honest, and do not distinguish permissibility with **beliefs** or **knowledge** ⚠

## Streams (cont'd)

- An infinite sequence of reports  $\rho$  is a **stream** for a world  $W$  if

$$\langle i, c, \varphi \rangle \in \rho \iff W, c \models P_i \varphi$$

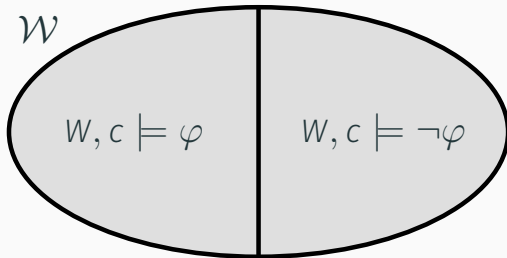
### Example



$$\rho = (\langle d, a, p \vee q \rangle, \langle d, a, p \rangle, \langle d, a, q \rangle, \\ \langle d, a, \neg p \rangle, \langle d, a, \neg q \rangle, \langle d, a, p \wedge q \rangle, \dots)$$

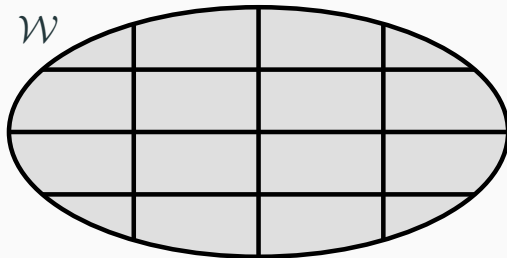
# Questions

- We want to design methods  $L$  which learn  $W$  when fed a stream  $\rho$
- Finding  $W$  exactly is too much to ask
- A **question**  $Q$  is a partition of  $\mathcal{W}$ 
  - $Q_{\varphi,c}$ : does  $\varphi$  hold in case  $c$ ?



# Questions

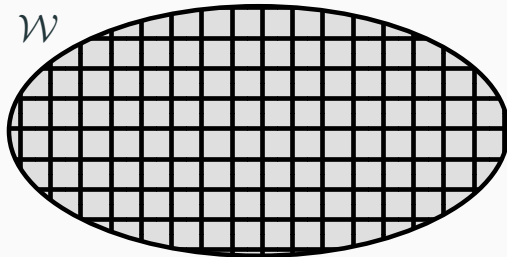
- We want to design methods  $L$  which learn  $W$  when fed a stream  $\rho$
- Finding  $W$  exactly is too much to ask
- A **question**  $Q$  is a partition of  $\mathcal{W}$ 
  - $Q_{\varphi,c}$ : does  $\varphi$  hold in case  $c$ ?
  - $Q_{\text{val}}$ : what are the correct valuations? (ignoring partitions)





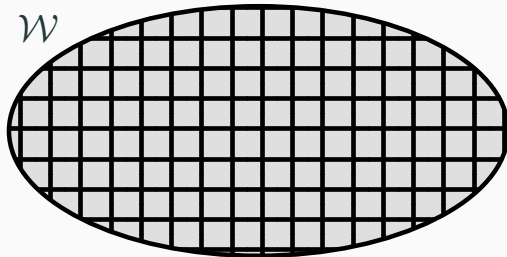
# Questions

- We want to design methods  $L$  which learn  $W$  when fed a stream  $\rho$
- Finding  $W$  exactly is too much to ask
- A **question**  $Q$  is a partition of  $\mathcal{W}$ 
  - $Q_{\varphi, c}$ : does  $\varphi$  hold in case  $c$ ?
  - $Q_{\text{val}}$ : what are the correct valuations? (ignoring partitions)
  - $Q_{\perp} = \{\{W\} \mid W \in \mathcal{W}\}$ : what is the actual world?



# Questions

- We want to design methods  $L$  which learn  $W$  when fed a stream  $\rho$
- Finding  $W$  exactly is too much to ask
- A **question**  $Q$  is a partition of  $\mathcal{W}$ 
  - $Q_{\varphi, c}$ : does  $\varphi$  hold in case  $c$ ?
  - $Q_{\text{val}}$ : what are the correct valuations? (ignoring partitions)
  - $Q_{\perp} = \{\{W\} \mid W \in \mathcal{W}\}$ : what is the actual world?
  - $Q[W]$  is the **correct answer** at  $W$



## Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

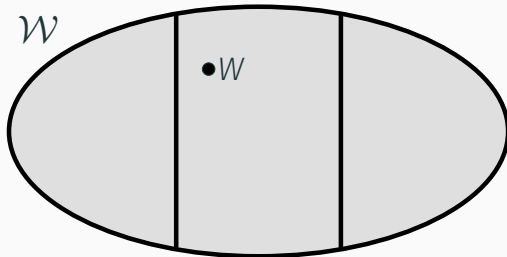
- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

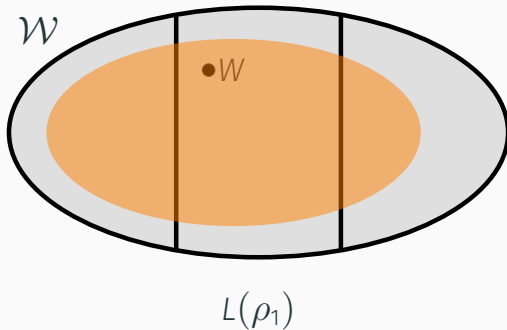


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

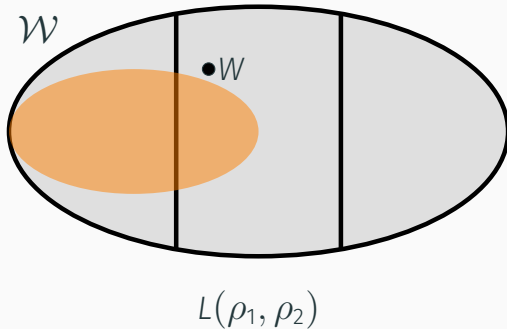


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

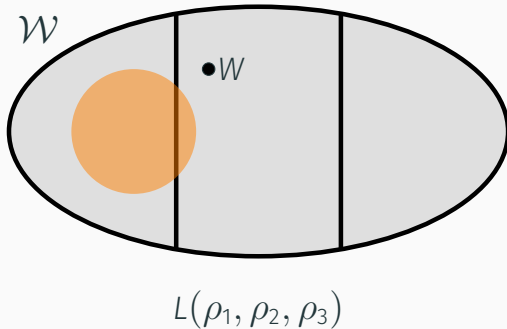


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

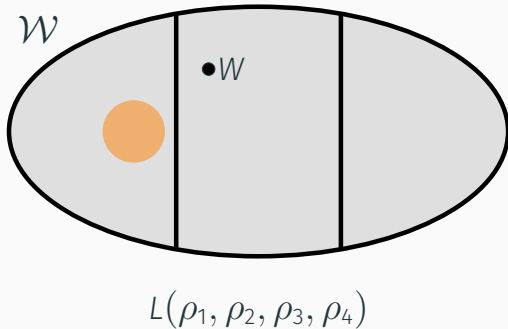


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$



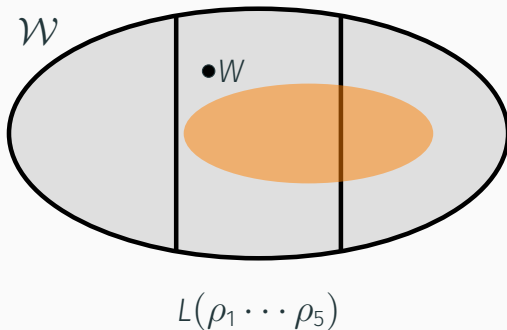


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

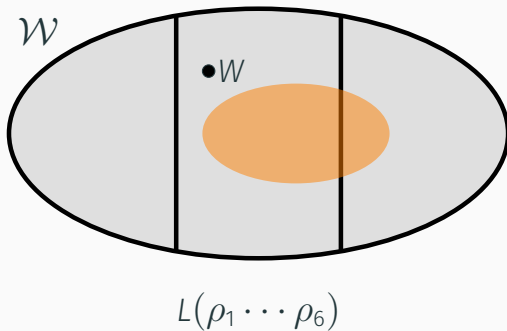


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

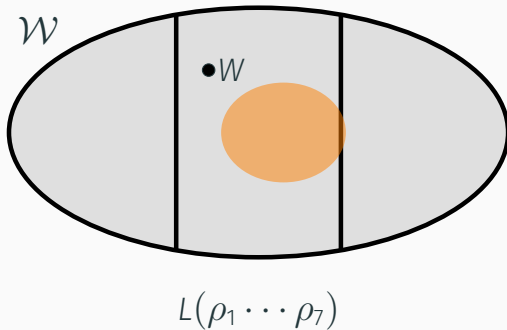


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

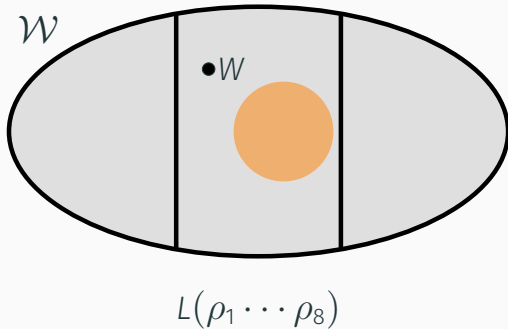


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

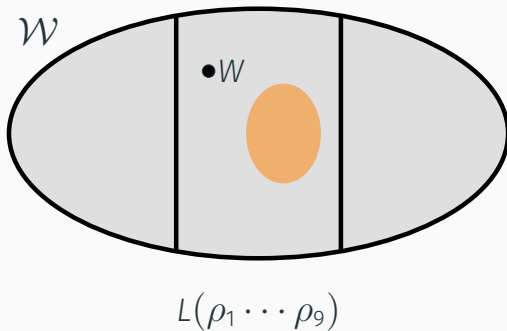


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$

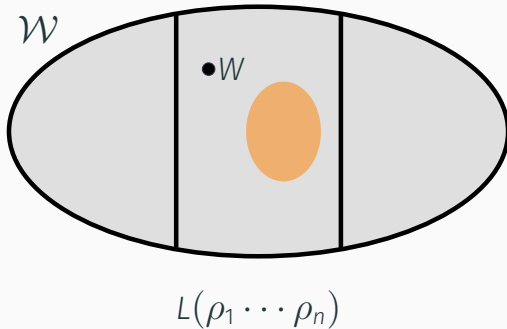


# Solvability

- $L$  **solves**  $Q$  if given any stream,  $L$  eventually finds the correct answer

$$\forall W, \forall \rho \text{ a stream for } W, \exists n \text{ s.t. } \forall m \geq n, L(\rho_1 \cdots \rho_m) \subseteq Q[W]$$

- $Q$  is **solvable** if there is a consistent method  $L$  which solves  $Q$



What can be learned?

---

## Solvable questions

- Which questions are solvable?
- It turns out there is a question  $Q^*$  which is the unique **hardest solvable question**

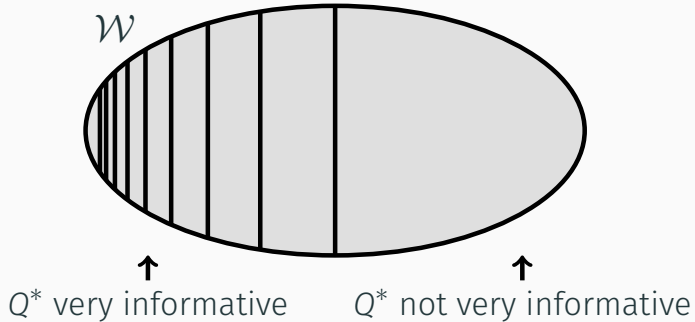
$$W \sim^* W' \iff \forall i \in \mathcal{S}, c \in \mathcal{C}, \Pi_i^W[v_c^W] = \Pi_i^{W'}[v_c^{W'}]$$

- Equivalently,  $W$  and  $W'$  have **exactly the same streams**
- $Q_{\varphi, c}$  is only solvable when  $\varphi$  is a tautology or contradiction **✗**
- $Q_{\text{val}}, Q_{\perp}$  not solvable **✗**
- **Problem:** if source have no expertise at all, *all* reports are permissible.  
True valuations don't matter!



## Solvable questions (cont'd)

- **Solution:** investigate what  $Q^*[W]$  tells us about  $W$



- A property of  $W$  is **learnable** if all  $W' \in Q^*[W]$  share the same property
  - Any method solving  $Q^*$  eventually finds it

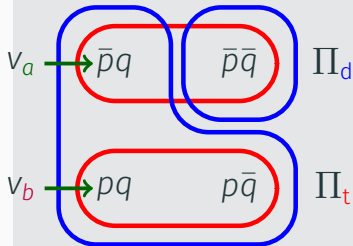
# What can be learned?

## Theorem

The true  $c$ -valuation is learnable at  $W$  iff there is a set  $\Gamma$  st

1.  $W, c \models \Gamma$
2.  $\text{Cn}(\Gamma)$  is a maximally consistent set
3. For all  $\varphi \in \Gamma$ , there is  $i \in \mathcal{S}$  such that  $W \models E_i \varphi$

## Example



For  $a$ , take  $\Gamma = \{p \vee q, \neg p\}$ :

$W, a \models \Gamma$ ,  $\text{Cn}(\Gamma) = \text{Cn}(\neg p \wedge q)$

$W \models E_d(p \vee q)$ ,  $W \models E_t \neg p$ ,



For  $b$  there is no such  $\Gamma$ !



- Similar result for partitions (omitted)

## Truth-tracking methods

---

## A characterisation of truth-tracking

- We have so far only looked solvable questions
- Which methods actually solve them?
- $L$  is **truth-tracking** if it solves all solvable questions
  - Equivalently,  $L$  solves  $Q^*$
- We characterise truth-tracking axiomatically, given three basic properties:
  - **Equivalence:** if  $\sigma \equiv \delta$  then  $L(\sigma) = L(\delta)$
  - **Repetition:**  $L(\sigma \cdots \sigma) = L(\sigma)$
  - **Permissibility:** if  $W \in L(\sigma)$  then  $W, c \models P_i\varphi$  for all  $\langle i, c, \varphi \rangle \in \sigma$

## A characterisation of truth-tracking (cont'd)

- Let  $T_\sigma$  be the set of worlds  $W$  such that, for all  $\langle i, c, \varphi \rangle$ :

$$W, c \models P_i\varphi \iff \exists\psi \equiv \varphi \text{ s.t. } \langle i, c, \psi \rangle \in \sigma$$

- i.e.  $\sigma$  contains all permissible reports, up to logical equivalence
- Write  $U, c \models \varphi$  iff  $W, c \models \varphi$  for all  $W \in U$
- **Credulity:** if  $T_\sigma, c \not\models P_i\varphi$  then  $L(\sigma), c \models \neg P_i\varphi$

### Theorem

For a method  $L$  satisfying *Equivalence*, *Repetition* and *Permissibility*,

$$\text{Truth-tracking} \iff \text{Credulity}$$

- **Credulity:** if  $T_\sigma, c \not\models P_i\varphi$  then  $L(\sigma), c \models \neg P_i\varphi$
- More expertise means **fewer permissible reports**
- **Credulity** is a principle of **maximal trust**
  - Whenever consistent with  $T_\sigma$ , we should trust  $i$  to have expertise to rule out  $\varphi$
  - Since all permissible reports eventually received, mistaken trust can be retracted
- **Consequence:** truth-tracking is not possible deductively; **inductive reasoning** is required
- Stronger property in terms of expertise directly: if  $T_\sigma \not\models \neg E_i\varphi$  then  $L(\sigma) \models E_i\varphi$

## Conclusion

---

- **Summary:**
  - Developed a logical framework to reason about expertise and permissible reports
  - Expressed a learning problem in this setting
  - Characterised conditions under which information can be learned
  - Axiomatically characterised truth-tracking learning methods
- **Future work:**
  - Assumptions on streams are very strong! Can these be lifted?
  - Everything is finite. What results carry over to the infinite case?
  - Bridge with probabilistic reasoning?



## An example method

- **Intuition:** express credulity with a prior **plausibility ordering** over worlds
- Conjecture the **maximally plausible worlds** consistent with permissibility statements
- E.g. using the number of partition cells as a measure of expertise:

$$L(\sigma) = \operatorname{argmax}_{W \in X_\sigma} \sum_{i \in \mathcal{S}} |\Pi_i^W|$$

where  $X_\sigma = \{W \mid \forall \langle i, c, \varphi \rangle \in \sigma, W, c \models P_i \varphi\}$

- This method *is* truth-tracking!