# Generalisation in Lifelong RL via Logical Composition

Geraud Nangue Tasse, Steven James, Benjamin Rosman

University of the Witwatersrand

**ICLR 2022**

# Motivation

- In general:
  - **Lifelong** agents that **reuse** past skills
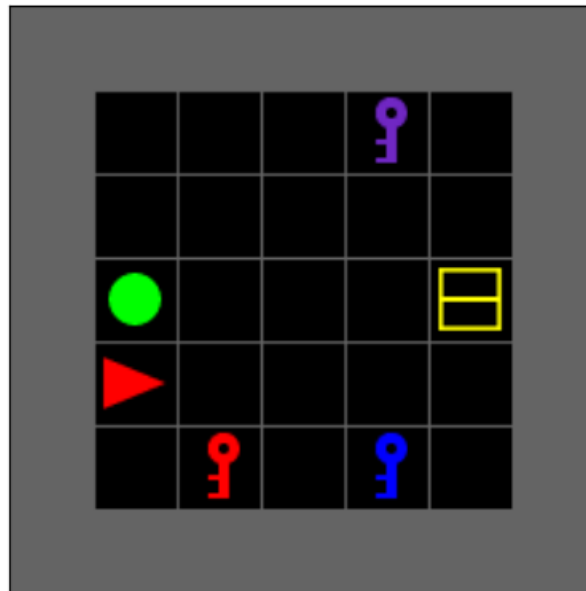
# Motivation

- In general:
  - **Lifelong** agents that **reuse** past skills

- In particular:
  - Are new tasks **expressible** in terms of learned ones?
    - If yes, **zero-shot learning**?
    - If no, **few-shot learning**?

# Motivation

- In general:
    - **Lifelong** agents that **reuse** past skills

- In particular:
    - Are new tasks **expressible** in terms of learned ones?
        - If yes, **zero-shot learning**?
        - If no, **few-shot learning**?
    - How about lifelong **generalisation**?

# Problem Setting

- Pickup-object domain

# Lifelong RL with Composition

| Goals | 🔑(red) | 🔴 | ⬟(red) | 🔑(blue) | 🔵 | ⬟(blue) | 🔑(green) | 🟢 | ⬟(green) | 🔑(purple) | 🟣 | ⬟(purple) | 🔑(yellow) | 🟡 | ⬟(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

*Learned*

# Lifelong RL with Composition

| Goals | 🔑(red) | ●(red) | 🎁(red) | 🔑(blue) | ●(blue) | 🎁(blue) | 🔑(green) | ●(green) | 🎁(green) | 🔑(purple) | ●(purple) | 🎁(purple) | 🔑(yellow) | ●(yellow) | 🎁(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

*Pick up Yellow Keys*

# Lifelong RL with Composition

| Goals | 🔑(red) | ●(red) | 🛢(red) | 🔑(blue) | ●(blue) | 🛢(blue) | 🔑(green) | ●(green) | 🛢(green) | 🔑(purple) | ●(purple) | 🛢(purple) | 🔑(yellow) | ●(yellow) | 🛢(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

*SOPGOL*

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🔴 | 🔑 | 🔵 | 🔵 | 🔑 | 🟢 | 🟢 | 🔑 | 🟣 | 🟣 | 🔑 | 🟡 | 🟡 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

Learned

*SOPGOL*

# Lifelong RL with Composition

| Goals | 🔑(red) | ●(red) | ▤(red) | 🔑(blue) | ●(blue) | ▤(blue) | 🔑(green) | ●(green) | ▤(green) | 🔑(purple) | ●(purple) | ▤(purple) | 🔑(yellow) | ●(yellow) | ▤(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

Learned

*SOPGOL*

Not green and not blue and yellow and key

➢ $T' := \neg\,🟩 \wedge \neg\,🟪 \wedge 🟨 \wedge 🔑$

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🔺 | 🔑 | 🔵 | 🔺 | 🔑 | 🟢 | 🔺 | 🔑 | 🟣 | 🔺 | 🔑 | 🟡 | 🔺 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

*Learned*

## *SOPGOL*

*Not green and not blue and yellow and key*

➢ $T' := \neg$🟩$\land \neg$🟪$\land$🟨$\land$🔑

➢ $Q_{SOP} = \neg Q^*(🟩) \land \neg Q^*(🟪) \land Q^*(🟨) \land Q^*(🔑)$

# Lifelong RL with Composition

| Goals | 🔑(red) | ⬤(red) | ☰(red) | 🔑(blue) | ⬤(blue) | ☰(blue) | 🔑(green) | ⬤(green) | ☰(green) | 🔑(purple) | ⬤(purple) | ☰(purple) | 🔑(yellow) | ⬤(yellow) | ☰(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

*Learned*

*SOPGOL*

*Not green and not blue and yellow and key*

➢ $T' := \neg$🟩$\land \neg$🟪$\land$🟨$\land$🔑

➢ $Q_{SOP} = \neg Q^*(🟩) \land \neg Q^*(🟪) \land Q^*(🟨) \land Q^*(🔑)$

➢ $T = T'$ ? **Yes!**

# Lifelong RL with Composition

| Goals | 🗝 | 🔴 | 🟥 | 🗝 | 🔵 | 🟦 | 🗝 | 🟢 | 🟩 | 🗝 | 🟣 | 🟪 | 🗝 | 🟡 | 🟨 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🗝 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

*Learned*

### SOPGOL

*Not green and not blue and yellow and key*

➤ $T' := \neg\,🟩 \wedge \neg\,🟪 \wedge 🟨 \wedge 🗝$

➤ $Q_{SOP} = \neg Q^*(🟩) \wedge \neg Q^*(🟪) \wedge Q^*(🟨) \wedge Q^*(🗝)$

➤ $T = T'$ ? **Yes!**

*Reduces the RL problem to just Supervised learning*

# Lifelong RL with Composition

| Goals | 🔑(red) | ●(red) | ▤(red) | 🔑(blue) | ●(blue) | ▤(blue) | 🔑(green) | ●(green) | ▤(green) | 🔑(purple) | ●(purple) | ▤(purple) | 🔑(yellow) | ●(yellow) | ▤(yellow) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ■(green) | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| ■(purple) | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ■(yellow) | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

*Pick up Yellow Boxes*

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🔴 | 🔑 | 🔵 | 🔵 | 🔑 | 🟢 | 🟢 | 🔑 | 🟣 | 🟣 | 🔑 | 🟡 | 🟡 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

*Learned*

*SOPGOL*

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🟥 | 🔑 | 🔵 | 🟦 | 🔑 | 🟢 | 🟩 | 🔑 | 🟣 | 🟪 | 🔑 | 🟡 | 🟨 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

Learned

*SOPGOL*

Not green and not blue and yellow and not key

> $T' \coloneqq \neg\,\blacksquare \wedge \neg\,\blacksquare \wedge \blacksquare \wedge \neg\,🔑$

> $Q_{SOP} = \neg Q^*(\blacksquare) \wedge \neg Q^*(\blacksquare) \wedge Q^*(\blacksquare) \wedge \neg Q^*(🔑)$

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🟥 | 🔑 | 🔵 | 🟦 | 🔑 | 🟢 | 🟩 | 🔑 | 🟣 | 🟪 | 🔑 | 🟡 | 🟨 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

*Learned*

*SOPGOL*

*Not green and not blue and yellow and not key*

➢ $T' := \neg\,🟩 \wedge \neg\,🟪 \wedge 🟨 \wedge \neg\,🔑$

➢ $Q_{SOP} = \neg Q^*(🟩) \wedge \neg Q^*(🟪) \wedge Q^*(🟨) \wedge \neg Q^*(🔑)$

➢ $T = T'$ ? **No!**

# Lifelong RL with Composition

| Goals | 🔑 | 🔴 | 🟥 | 🔑 | 🔵 | 🟦 | 🔑 | 🟢 | 🟩 | 🔑 | 🟣 | 🟪 | 🔑 | 🟡 | 🟨 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 🟩 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟪 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 🟨 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 🔑 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

*Learned*

*SOPGOL*

*Not green and not blue and yellow and not key*

➢ $T' := \neg\,🟩 \wedge \neg\,🟪 \wedge 🟨 \wedge \neg\,🔑$

➢ $Q_{SOP} = \neg Q^*(🟩) \wedge \neg Q^*(🟪) \wedge Q^*(🟨) \wedge \neg Q^*(🔑)$

➢ $T = T'$ ? **No!**

➢ Learn new $Q$ with goal-oriented learning **(using $Q_{SOP}$, to speed up training)**, then add to library

# Transfer After Pretraining

- Pretrained: , , , 

| Goals | 🔑 | ⬤ | ▤ | 🔑 | ⬤ | ▤ | 🔑 | ⬤ | ▤ | 🔑 | ⬤ | ▤ | 🔑 | ⬤ | ▤ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T_2$ | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |

*Pick up red/blue/green boxes, or green/purple/yellow keys, or blue balls*

# Transfer After Pretraining

- Pretrained: 🟩, 🟪, 🟨, 🗝

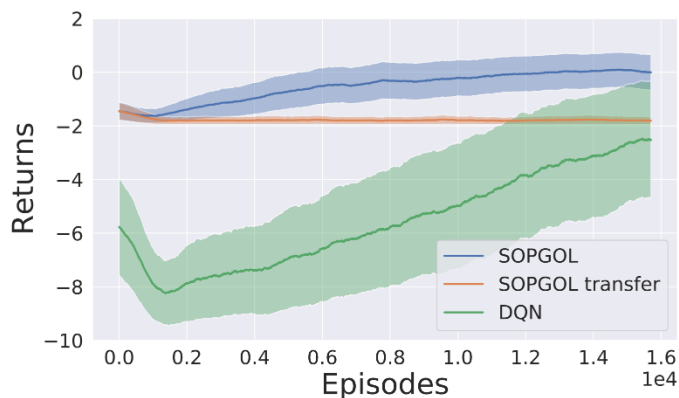| Goals | 🗝 | ⬤ | 🟥 | 🗝 | ⬤ | 🟦 | 🗝 | ⬤ | 🟩 | 🗝 | ⬤ | 🟪 | 🗝 | ⬤ | 🟨 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $T_2$ | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |



*SOPGOL > DQN*

# Transfer After Pretraining

- Pretrained: 🟩, 🟪, 🟨, 🗝

| Goals | 🔑 | 🔴 | 🟥 | 🔑 | 🔵 | 🟦 | 🔑 | 🟢 | 🟩 | 🔑 | 🟣 | 🟪 | 🔑 | 🟡 | 🟨 |
|-------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| $T_2$ | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |



*SOPGOL > DQN*

Inferred behaviour: $(🟩 \wedge \neg 🟪 \wedge \neg 🟨) \vee (\neg 🟩 \wedge 🟪 \wedge \neg 🟨 \wedge 🗝) \vee (\neg 🟩 \wedge \neg 🟪 \wedge 🟨 \wedge 🗝) \vee (\neg 🟪 \wedge \neg 🟨 \wedge \neg 🗝)$

# Transfer After Pretraining

- Pretrained: 🟩, 🟪, 🟨, 🔑



SOPGOL > DQN

Learned representation of new task

Inferred behaviour:
$$(\square \wedge \neg \square \wedge \neg \square) \vee (\neg \square \wedge \square \wedge \neg \square \wedge \text{🔑}) \vee (\neg \square \wedge \neg \square \wedge \square \wedge \text{🔑}) \vee (\neg \square \wedge \neg \square \wedge \neg \text{🔑})$$

# Lifelong RL with Composition

*SOPGOL*

**Theorem**:    $\log(|goals|) \leq \lim_{t \to \infty}(|\boldsymbol{skills}|) \leq |goals|$

Note:    $|\text{tasks}| = 2^{|goals|}$

# Generalisation with Lifelong Transfer

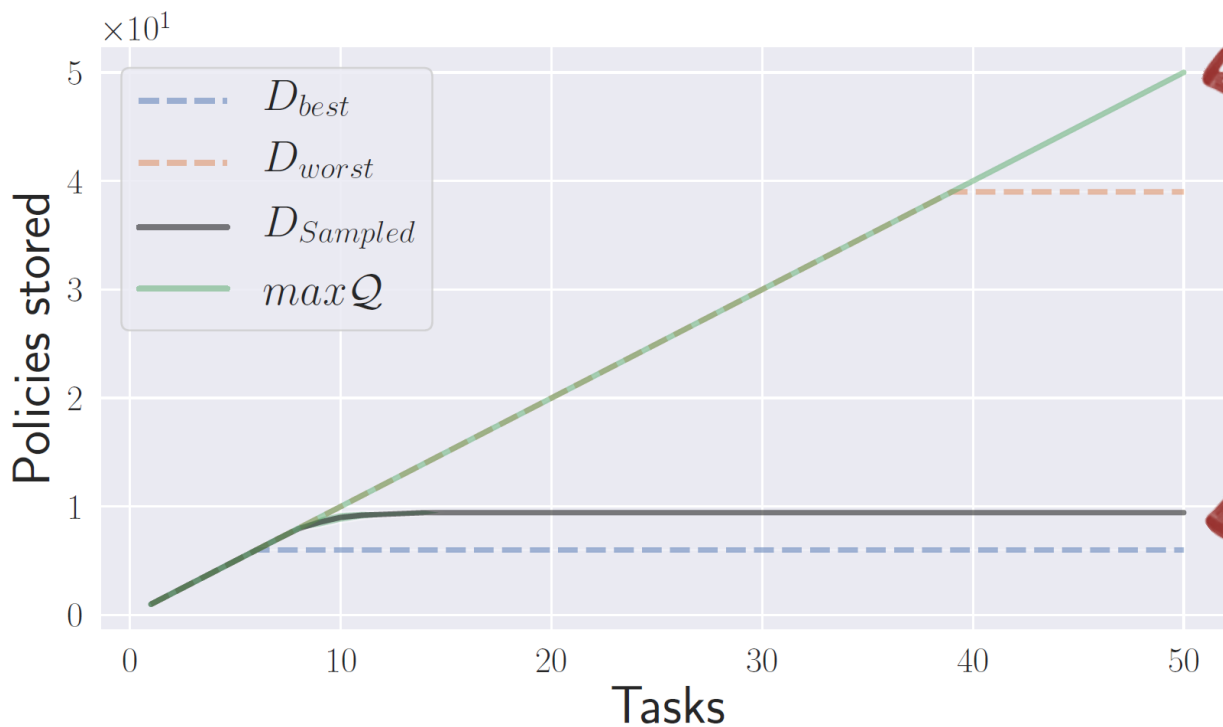- Four-rooms domain: 40 goals, $2^{40}$ ~ 1 trillion tasks

# Conclusion

- We leverage logical composition for **fast transfer** between tasks.

- Leads to **quick generalisation** over any task distribution.
  - we have a **logarithmic upper bound** on the number of tasks that needs to be learned and stored.

- Leads to both **interpretable** and **sample-efficient** lifelong RL.