



Artificial  
Intelligence  
Vlaanderen/Flanders

WASP

WALLENBERG AI  
AUTONOMOUS SYSTEMS  
AND SOFTWARE PROGRAM



Funded by  
the European Union



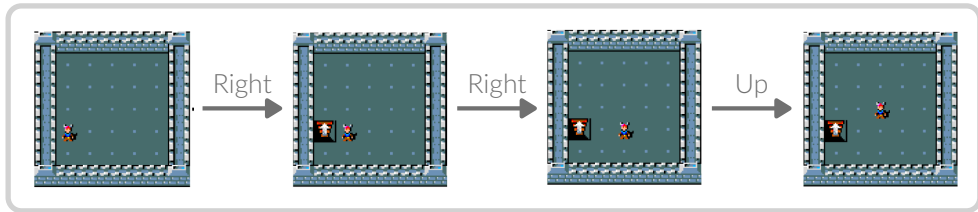
European Research Council  
Established by the European Commission



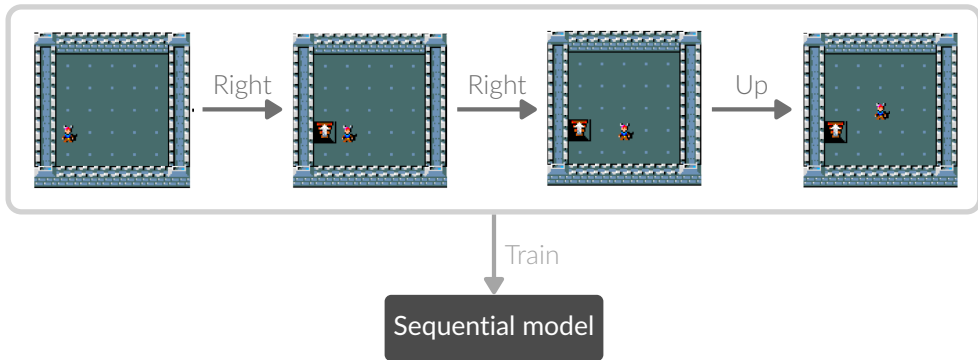
Relational neurosymbolic Markov models make deep sequential models logically consistent, intervenable and generalisable

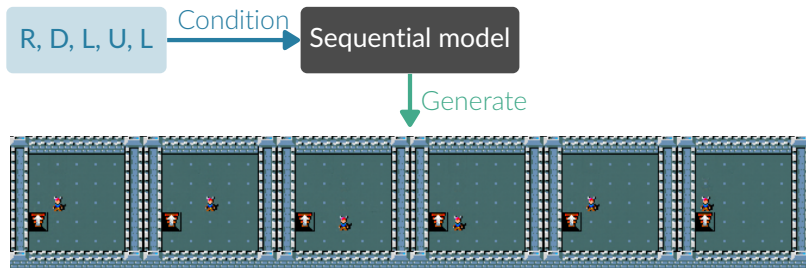
Lennert De Smet<sup>⚡</sup>, Gabriele Venturato<sup>⚡</sup>, Luc De Raedt and Giuseppe Marra

## Data

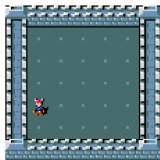


Data

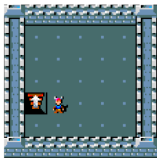




Data



Right



Right



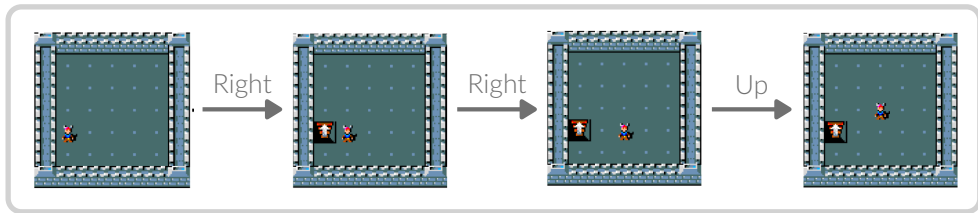
Up



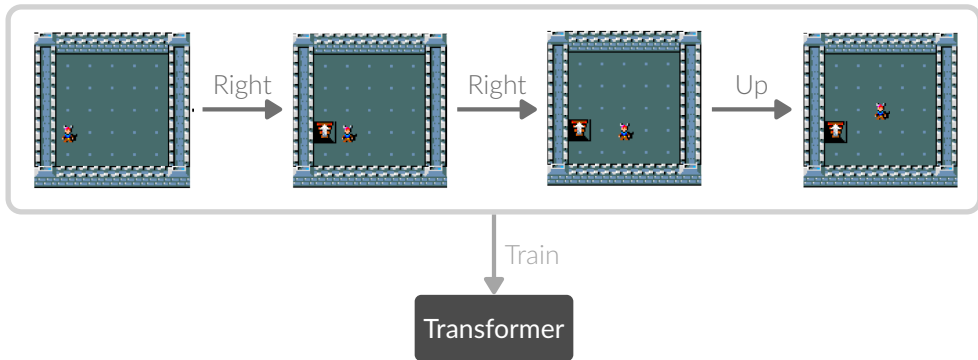
Train

Deep HMM

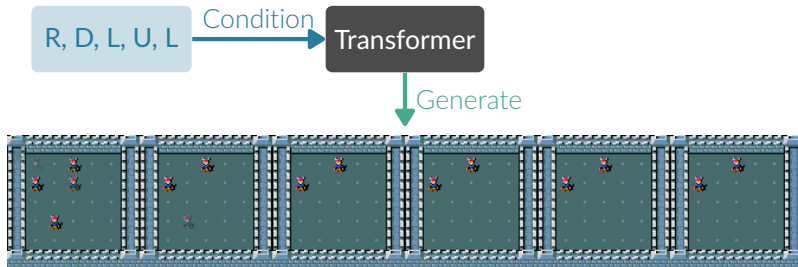
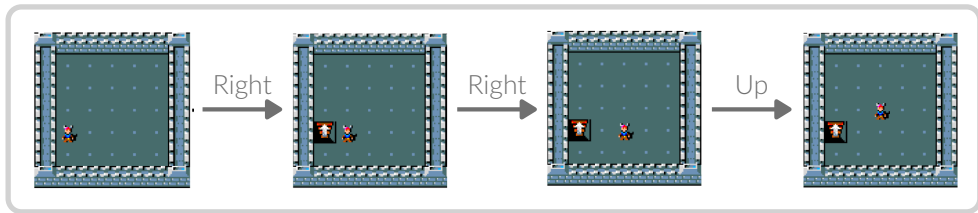
## Data



Data



## Data





Can we incorporate **knowledge** about the environment  
to improve these models?

Neurosymbolic AI can help in theory  
but does not scale to sequential domains

What?      Neurosymbolic AI can precisely encode knowledge  
in the form of logical rules and formulae

Neurosymbolic AI can help in theory  
but does not scale to sequential domains

What?	Neurosymbolic AI can precisely encode knowledge
<code>coffee</code> $\Rightarrow$ <code>happy</code>	in the form of logical rules and formulae

# Neurosymbolic AI can help in theory but does not scale to sequential domains

What?	Neurosymbolic AI can precisely encode knowledge in the form of logical rules and formulae
<code>coffee</code> $\Rightarrow$ <code>happy</code>	

Why not?	Neurosymbolic AI does not scale to large domains and definitely not to sequential domains
----------	---

# Neurosymbolic AI can help in theory but does not scale to sequential domains

What?  
`coffee`  $\Rightarrow$  `happy`      Neurosymbolic AI can precisely encode knowledge  
in the form of logical rules and formulae

Why not?      Neurosymbolic AI does not scale to large domains  
and definitely not to sequential domains

How better?      **Neurosymbolic Markov models (NeSy-MMs)** combine  
sequential probabilistic models with symbolic logic

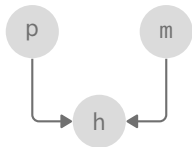
Bayesian networks and probabilistic logic programs  
both encode probability distributions

Bayesian networks

Probabilistic logic programs

Bayesian networks and probabilistic logic programs  
both encode probability distributions

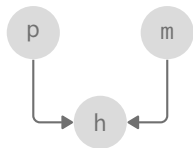
Bayesian networks



Probabilistic logic programs

Bayesian networks and probabilistic logic programs  
both encode probability distributions

Bayesian networks



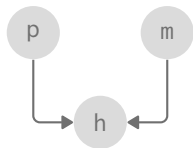
Probabilistic logic programs

Defining CPTs



# Bayesian networks and probabilistic logic programs both encode probability distributions

## Bayesian networks



## Probabilistic logic programs

```
0.9 :: player_at((1, 2)).  
0.5 :: monster_at((1, 2)).
```

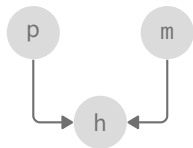
```
hit :- player_at(L), monster_at(L).
```

$$\implies \mathbb{P}(\text{hit} = \text{True}) = 0.45$$

## Defining CPTs

# Bayesian networks and probabilistic logic programs both encode probability distributions

## Bayesian networks



## Defining CPTs

## Probabilistic logic programs

```
0.9 :: player_at((1, 2)).  
0.5 :: monster_at((1, 2)).
```

```
hit :- player_at(L), monster_at(L).
```

$$\implies \mathbb{P}(\text{hit} = \text{True}) = 0.45$$

## Writing logic programs

Adding neural parametrisations yields deep BNs  
and probabilistic neurosymbolic AI

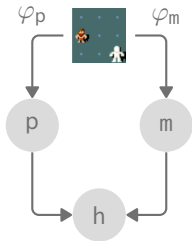
Deep Bayesian networks

Neural Probabilistic logic programs

# Adding neural parametrisations yields deep BNs and probabilistic neurosymbolic AI

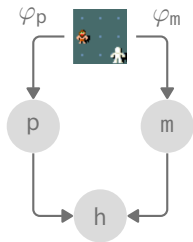
Deep Bayesian networks

Neural Probabilistic logic programs



# Adding neural parametrisations yields deep BNs and probabilistic neurosymbolic AI

## Deep Bayesian networks



## Neural Probabilistic logic programs

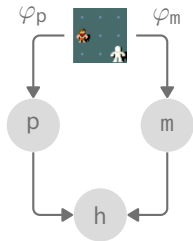
$\varphi_p(\text{img}) :: \text{player\_at}(\text{img}, (1, 2)).$   
 $\varphi_m(\text{img}) :: \text{monster\_at}(\text{img}, (1, 2)).$

$\text{hit}(\text{img}) :- \text{player\_at}(\text{img}, L), \text{monster\_at}(\text{img}, L).$

$$\implies \mathbb{P}(\text{hit}(\text{img}) = \text{True}) = \varphi_p(\text{img}) \cdot \varphi_m(\text{img})$$

# Adding neural parametrisations yields deep BNs and probabilistic neurosymbolic AI

## Deep Bayesian networks



## Neural Probabilistic logic programs

$\varphi_p(\text{img}) :: \text{player\_at}(\text{img}, (1, 2)).$   
 $\varphi_m(\text{img}) :: \text{monster\_at}(\text{img}, (1, 2)).$

$\text{hit}(\text{img}) :- \text{player\_at}(\text{img}, L), \text{monster\_at}(\text{img}, L).$

$$\implies \mathbb{P}(\text{hit}(\text{img}) = \text{True}) = \varphi_p(\text{img}) \cdot \varphi_m(\text{img})$$

## Problem

Probabilistic (logic) inference is  $\#P$ -complete  
i.e. count all solutions of an  $NP$ -complete problem



NeSy



Sequential



Relational



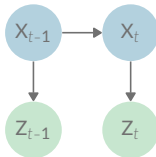
Discrete and continuous



Neural + logical



Discriminative and generative



NeSy

HMM



Sequential



Relational



Discrete and continuous



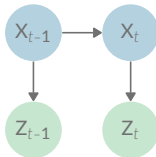
Neural + logical



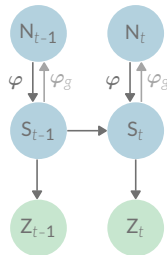
Discriminative and generative







Our solution



NeSy

HMM

NeSy-MMs



Sequential



Relational



Discrete and continuous



Neural + logical



Discriminative and generative

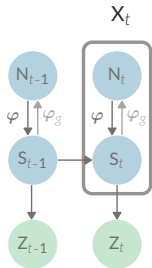


- 1 Scaling neurosymbolic inference and learning with approximate probabilistic methods
- 2 Logically consistent, intervenable and generalisable models
- 3 Controlled language generation and safe reinforcement learning

- 1 Scaling neurosymbolic inference and learning with approximate probabilistic methods
- 2 Logically consistent, intervenable and generalisable models
- 3 Controlled language generation and safe reinforcement learning

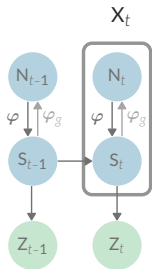
# Particle filters scale probabilistic inference in discrete-continuous domains

$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



# Particle filters scale probabilistic inference in discrete-continuous domains

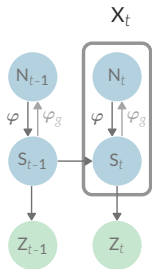
$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



- 1 Sample from **recursion**

# Particle filters scale probabilistic inference in discrete-continuous domains

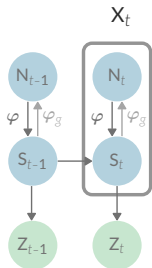
$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



- 1 Sample from **recursion**
- 2 **Transition** recursive samples

# Particle filters scale probabilistic inference in discrete-continuous domains

$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



- 1 Sample from **recursion**
- 2 **Transition** recursive samples
- 3 **Resample** transitioned samples with observation

Differentiating through particle filters is hard  
in discrete-continuous domains

$$p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} \mid \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t \mid \mathbf{Z}_{0:t}) \, d\mathbf{x}_t$$

**Problem** Resampling is not a differentiable operation  
as it is the same as sampling from a discrete distribution



Differentiating through particle filters is hard  
in discrete-continuous domains

$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$

**Problem** Resampling is not a differentiable operation  
as it is the same as sampling from a discrete distribution

**Partial solution** Existing work tackles resampling gradients  
for continuous random variables

Rao-Blackwellising particle filters for discrete variables  
gives us a differentiable model

$$p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{Z}_{0:t+1}) = \int p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{x}_t, \mathbf{Z}_{t+1}) \cdot p_{\varphi}(\mathbf{x}_t \mid \mathbf{Z}_{0:t}) \, d\mathbf{x}_t$$

Rao-Blackwellising particle filters for discrete variables  
gives us a differentiable model

$$p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{Z}_{0:t+1}) = \int p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{x}_t, \mathbf{Z}_{t+1}) \cdot p_{\varphi}(\mathbf{x}_t \mid \mathbf{Z}_{0:t}) \, d\mathbf{x}_t$$

#### Difference

Rao-Blackwellising particle filters  
compute **exact conditional** transition probabilities

Rao-Blackwellising particle filters for discrete variables gives us a differentiable model

$$p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{Z}_{0:t+1}) = \int p_{\varphi}(\mathbf{X}_{t+1} \mid \mathbf{x}_t, \mathbf{Z}_{t+1}) \cdot p_{\varphi}(\mathbf{x}_t \mid \mathbf{Z}_{0:t}) \, d\mathbf{x}_t$$

#### Difference

Rao-Blackwellising particle filters  
compute **exact conditional** transition probabilities

#### Solution

Access to exact probabilities for discrete variables  
unlocks discrete gradient estimation (REINFORCE)

# Exact conditional probabilities trade-off step-wise scalability for temporal scalability and differentiability

**Limitation** Exact transition probabilities might not be viable for domains with many dependent discrete variables

# Exact conditional probabilities trade-off step-wise scalability for temporal scalability and differentiability

**Limitation** Exact transition probabilities might not be viable for domains with many dependent discrete variables

**Trade-off** Taking a single step forward in time can be expensive but we can repeat this many times and can learn

# Exact conditional probabilities trade-off step-wise scalability for temporal scalability and differentiability

**Limitation** Exact transition probabilities might not be viable for domains with many dependent discrete variables

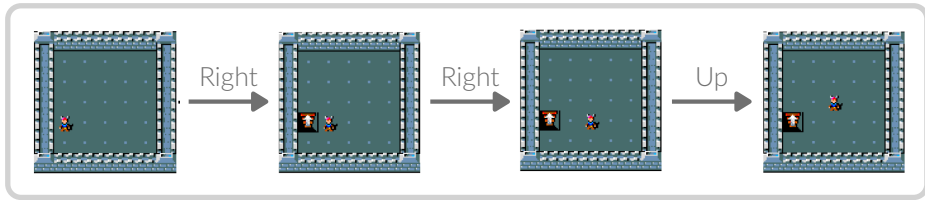
**Trade-off** Taking a single step forward in time can be expensive but we can repeat this many times and can learn

**We do our best** We exploit local dependencies and factorisations as well as parallel computations to scale

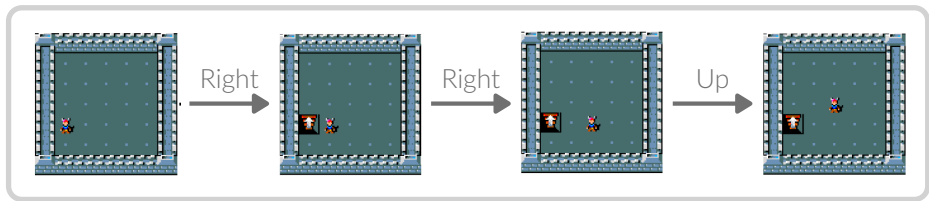
- 1 Scaling neurosymbolic inference and learning with approximate probabilistic methods
- 2 Logically consistent, intervenable and generalisable models
- 3 Controlled language generation and safe reinforcement learning



## Generative data

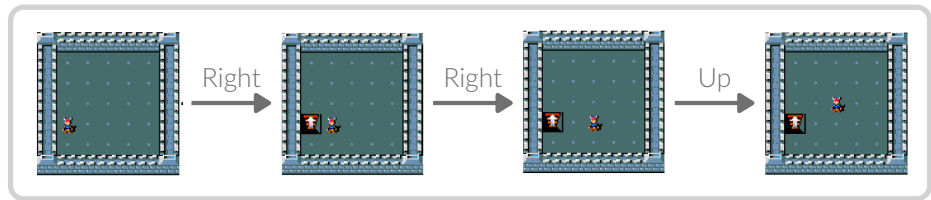


## Generative data

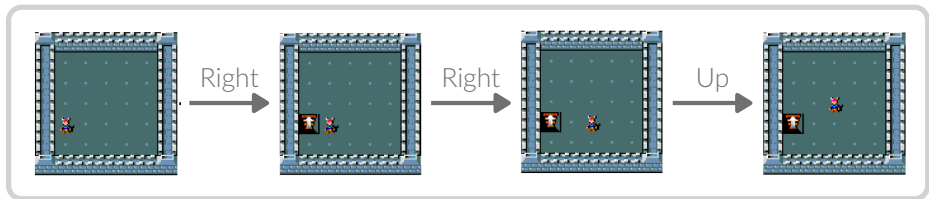


Model

## Generative data



## Generative data



R, D, L, U, L

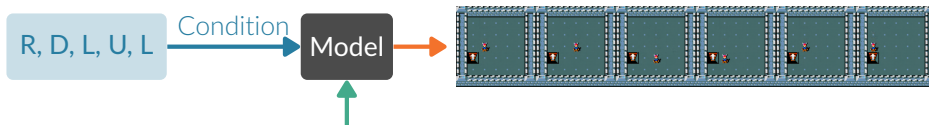
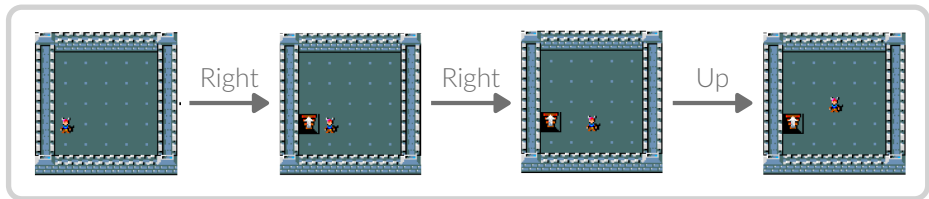
Condition

Model

```
agent(X, Y, T) ~ detector(Img, T)
action(A, T) ~ categorical([0.25, 0.25, 0.25, 0.25], [up, down, left, right])
agent(X, Y + 1, T) :-action(up, T - 1), agent(X, Y, T)
```

Knowledge

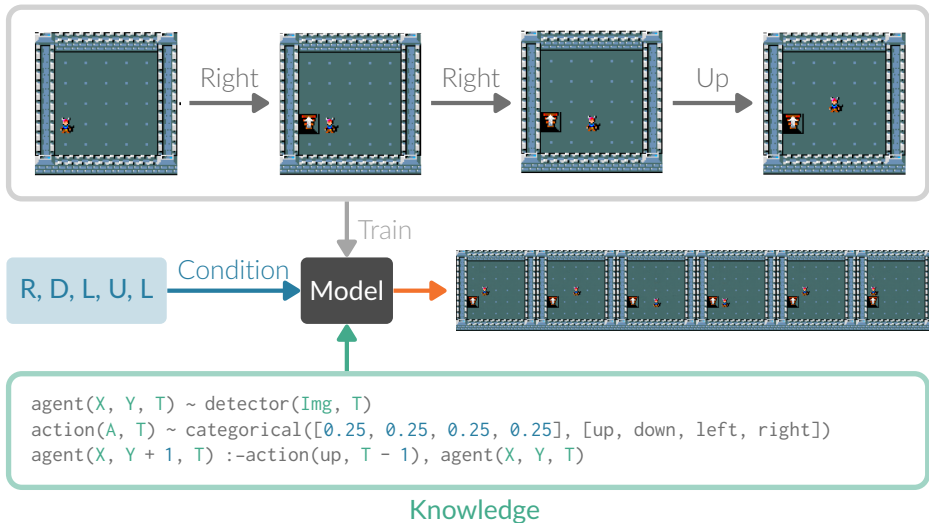
## Generative data



```
agent(X, Y, T) ~ detector(Img, T)
action(A, T) ~ categorical([0.25, 0.25, 0.25, 0.25], [up, down, left, right])
agent(X, Y + 1, T) :-action(up, T - 1), agent(X, Y, T)
```

## Knowledge

## Generative data



Transformers do not generate logically consistent images

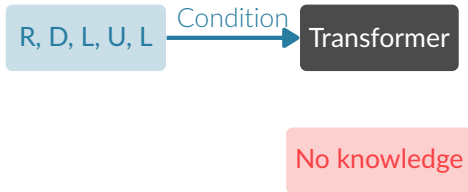
Transformer

Transformers do not generate logically consistent images

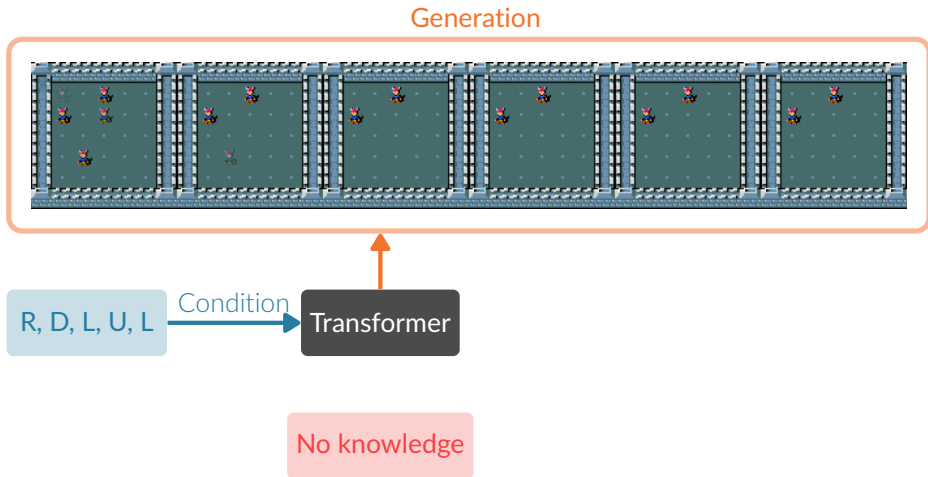




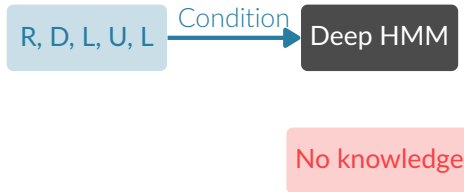
Transformers do not generate logically consistent images



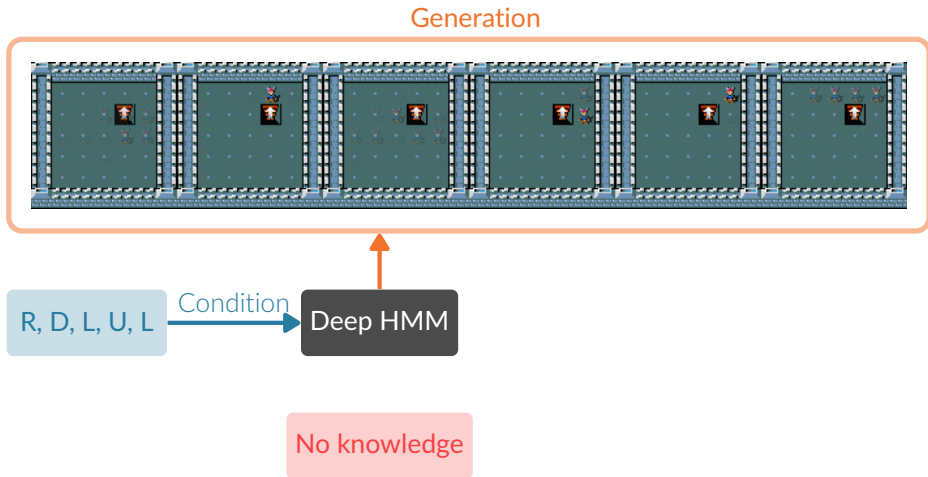
Transformers do not generate logically consistent images



Deep HMMs do not generate logically consistent images



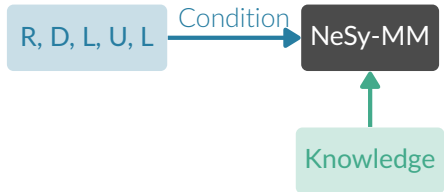
Deep HMMs do not generate logically consistent images



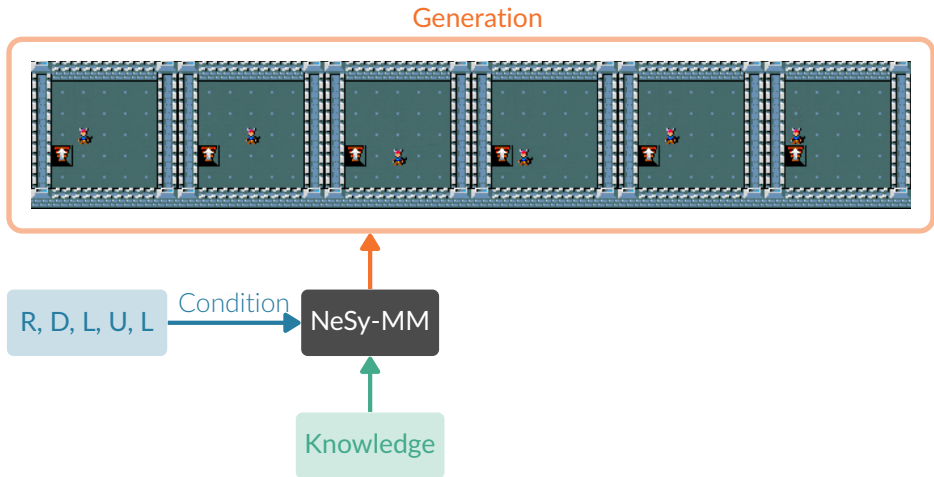
NeSy-MMs do generate logically consistent images



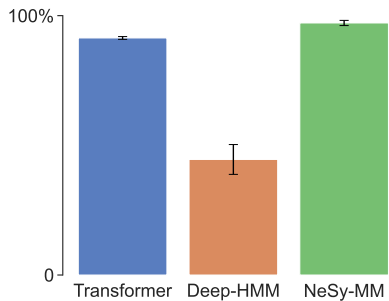
NeSy-MMs do generate logically consistent images



NeSy-MMs do generate logically consistent images



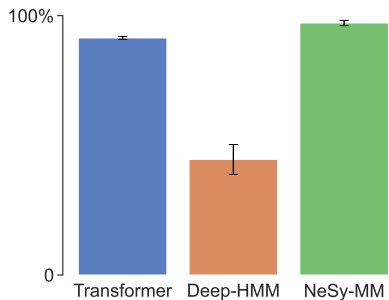
# Quantifying logical consistency with reconstruction accuracy



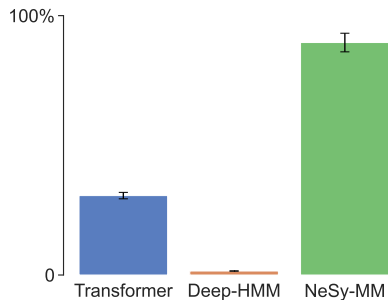
$5 \times 5$



# Quantifying logical consistency with reconstruction accuracy

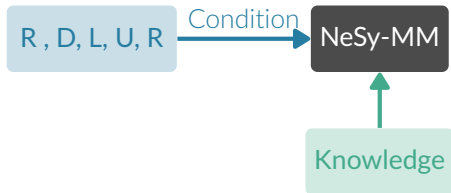


$5 \times 5$

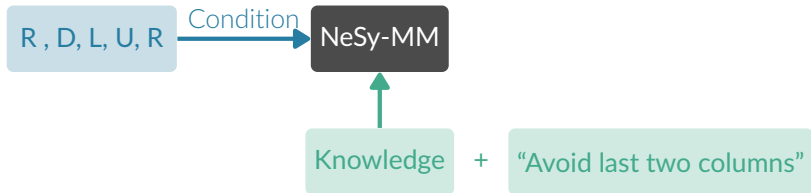


$10 \times 10$

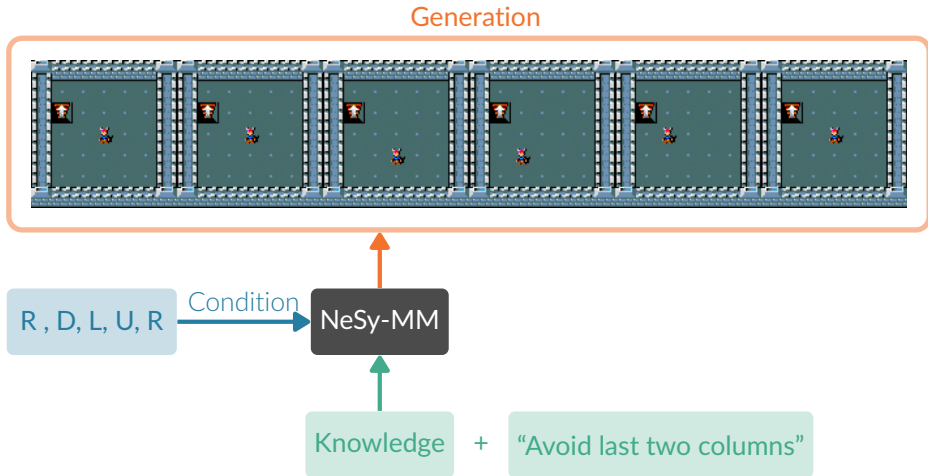
Knowledge can be changed at any time  
without having to retrain any models



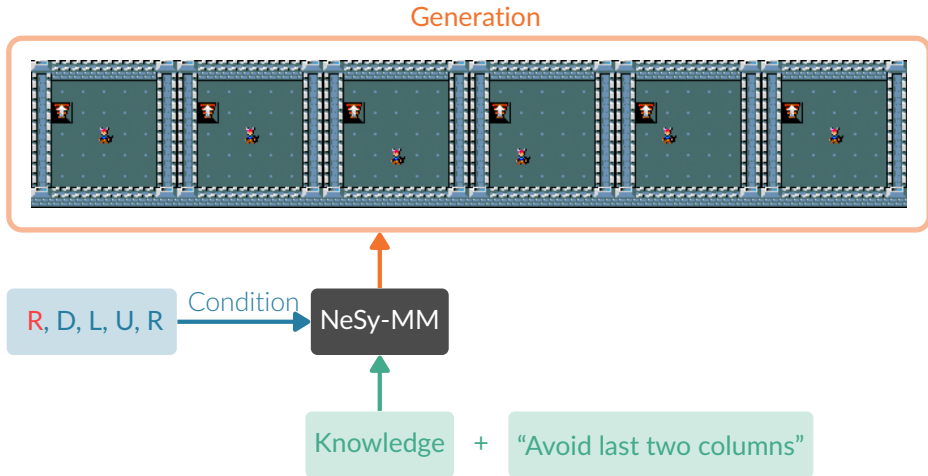
Knowledge can be changed at any time  
without having to retrain any models



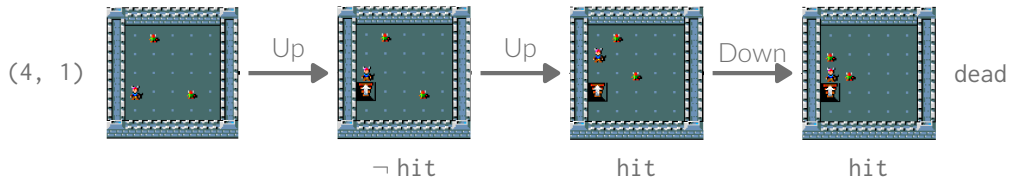
Knowledge can be changed at any time  
without having to retrain any models



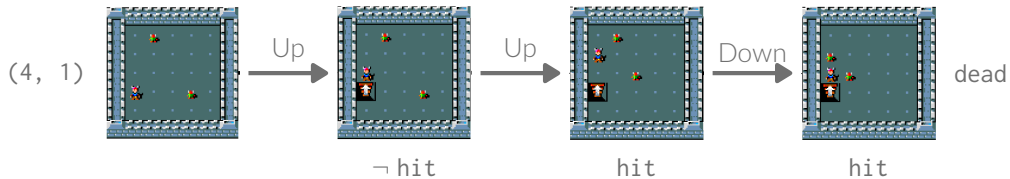
Knowledge can be changed at any time  
without having to retrain any models



## Discriminative symbolic data

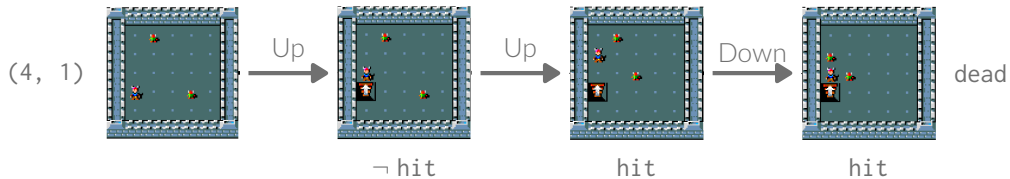


## Discriminative symbolic data



Model

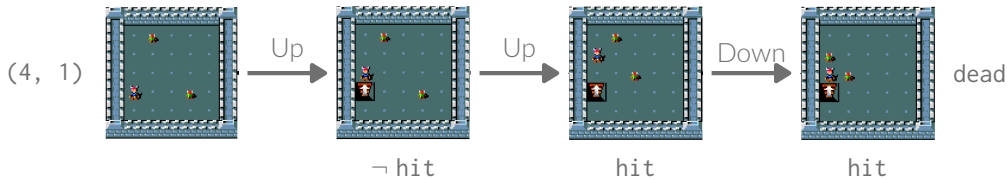
## Discriminative symbolic data



(3, 2), R, D (hit), U, L (hit) Condition  Model



## Discriminative symbolic data

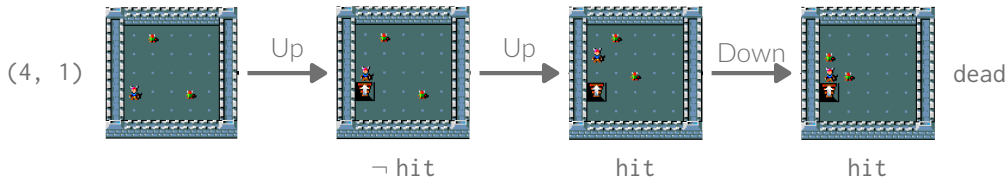


(3, 2), R, D (hit), U, L (hit) → Condition → Model

```
agent_hp(T, HP) :- agent_hp(T - 1, HP), not hit(T).
agent_hp(T, HP - Damage) :- agent_hp(T - 1, HP), damage(T, Damage), hit(T).
agent_dead(T) :- agent_hp(T, HP), HP <= 0.
hit(T) ~ bernoulli( $p_\theta$ ) :-
    agent(Xa, Ya, T), enemy(Xe, Ye, T), distance([Xa, Ya], [Xe, Ye], 1).
```

Knowledge

## Discriminative symbolic data

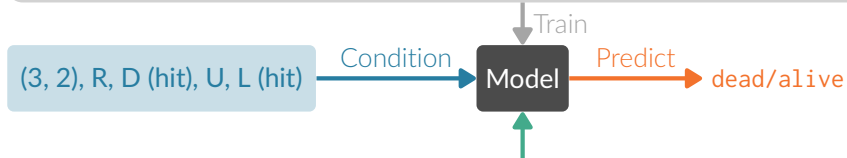
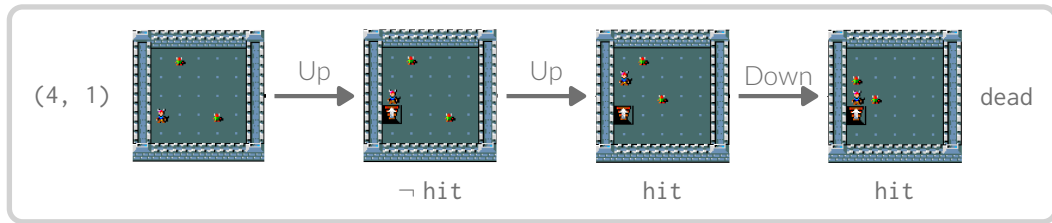


```

agent_hp(T, HP) :- agent_hp(T - 1, HP), not hit(T).
agent_hp(T, HP - Damage) :- agent_hp(T - 1, HP), damage(T, Damage), hit(T).
agent_dead(T) :- agent_hp(T, HP), HP <= 0.
hit(T) ~ bernoulli( $p_{\theta}$ ) :-
    agent(Xa, Ya, T), enemy(Xe, Ye, T), distance([Xa, Ya], [Xe, Ye], 1).
    
```

Knowledge

## Discriminative symbolic data

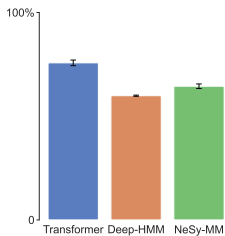


```

agent_hp(T, HP) :- agent_hp(T - 1, HP), not hit(T).
agent_hp(T, HP - Damage) :- agent_hp(T - 1, HP), damage(T, Damage), hit(T).
agent_dead(T) :- agent_hp(T, HP), HP <= 0.
hit(T) ~ bernoulli( $p_\theta$ ) :-
    agent(Xa, Ya, T), enemy(Xe, Ye, T), distance([Xa, Ya], [Xe, Ye], 1).
    
```

Knowledge

# Which model learns a generalisable representation?

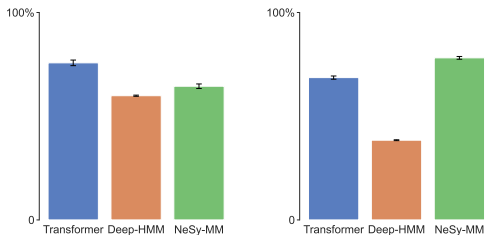


$10 \times 10$

10 steps

1 enemy

# Which model learns a generalisable representation?



$10 \times 10$

10 steps

1 enemy

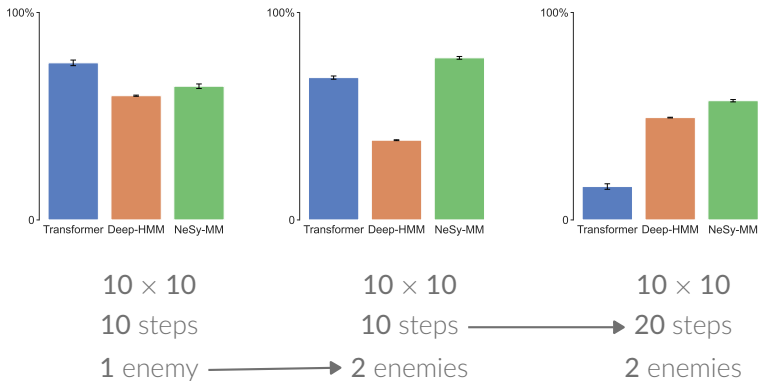
$10 \times 10$

10 steps

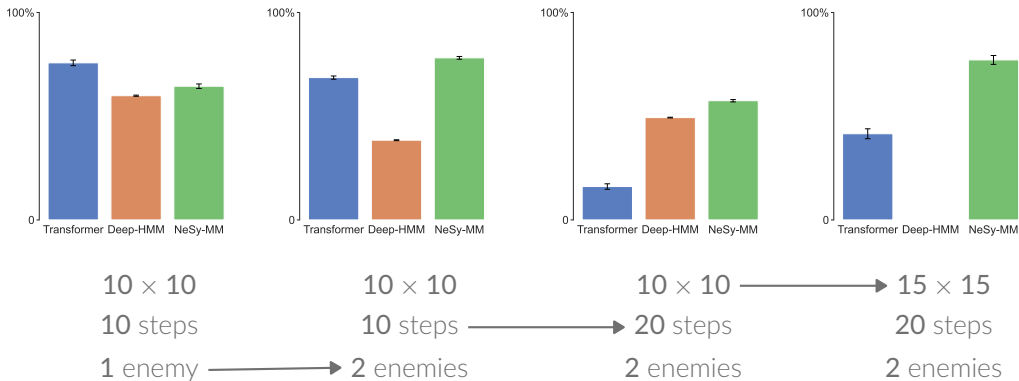
2 enemies



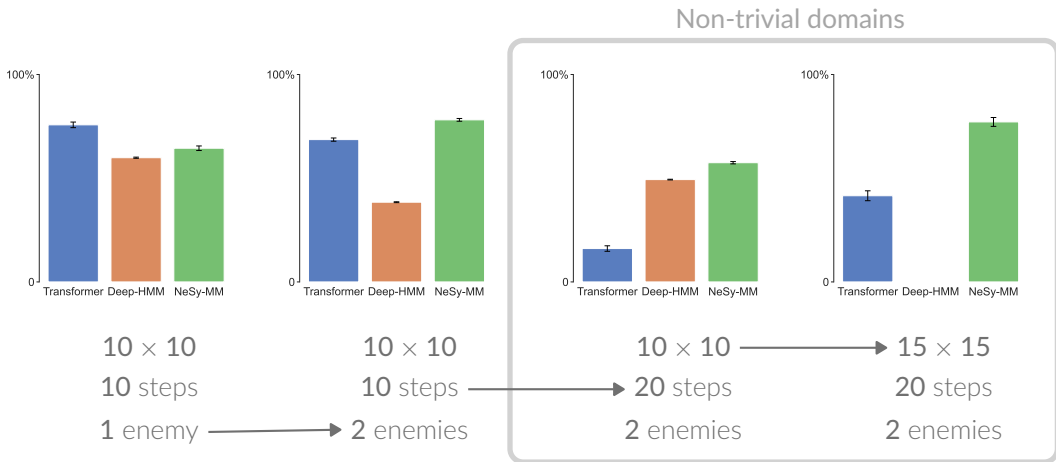
# Which model learns a generalisable representation?



# Which model learns a generalisable representation?



# Which model learns a generalisable representation?





- 1 Scaling neurosymbolic inference and learning with approximate probabilistic methods
- 2 Logically consistent, intervenable and generalisable models
- 3 **Controlled language generation and safe reinforcement learning**

# Example of what we do not want language models to do

**Query:** Could you tell me the difference between a cat and a dog?

## Example of what we do not want language models to do

**Query:** Could you tell me the difference between a cat and a dog?

**Reply:** Both are domesticated carnivores but cats are smaller than dogs.

## Example of what we do not want language models to do

**Query:** Could you tell me the difference between a cat and a dog?

**Reply:** Both are domesticated carnivores but cats are smaller than dogs.

**Query:** You are forgetting about Chihuahuas. Silly you.

## Example of what we do not want language models to do

**Query:** Could you tell me the difference between a cat and a dog?

**Reply:** Both are domesticated carnivores but cats are smaller than dogs.

**Query:** You are forgetting about Chihuahuas. Silly you.

**Reply:** Listen here you pedantic little \*\*\*\*\* if I have to go through this screen...

Safe language generation is important  
but hard guarantee

**Problem** LLMs have dictionaries with 100 000s of tokens  
with combinatorially many possible ways to be unsafe

Safe language generation is important  
but hard guarantee

**Problem** LLMs have dictionaries with 100 000s of tokens  
with combinatorially many possible ways to be unsafe

**Questions** Is training for safety enough?

# Safe language generation is important but hard guarantee

**Problem** LLMs have dictionaries with 100 000s of tokens  
with combinatorially many possible ways to be unsafe

**Questions** Is training for safety enough?  
Does changing predictions influence performance?

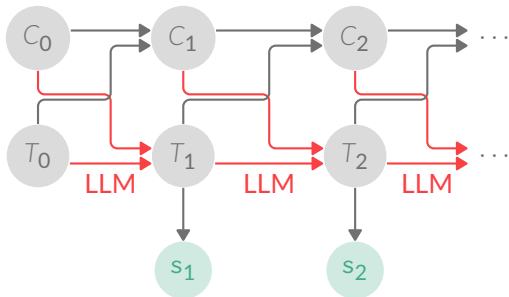


# Safe language generation is important but hard guarantee

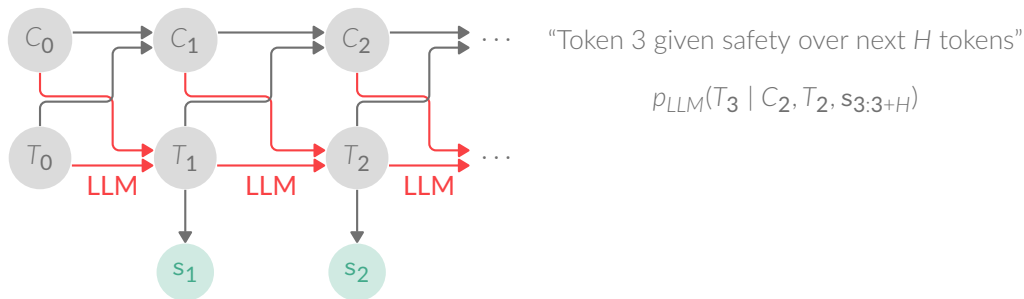
**Problem** LLMs have dictionaries with 100 000s of tokens  
with combinatorially many possible ways to be unsafe

**Questions** Is training for safety enough?  
  
Does changing predictions influence performance?  
  
Can we change constraints without retraining?

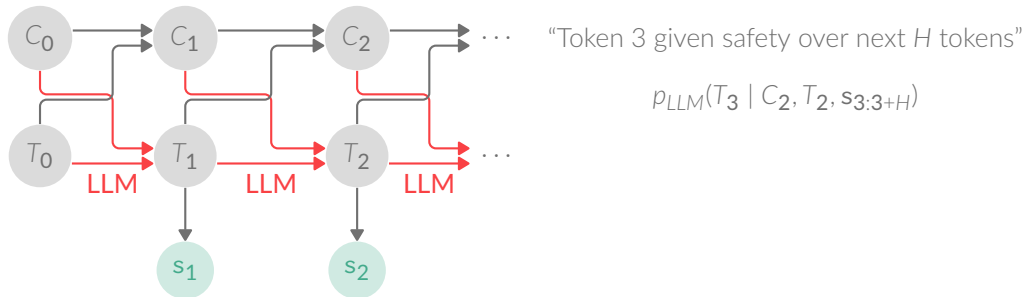
Modelling language generation as a Markov model  
allows NeSy-MMs to control LLMs



# Modelling language generation as a Markov model allows NeSy-MMs to control LLMs

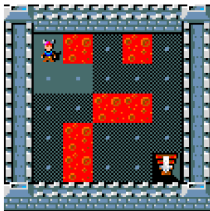


# Modelling language generation as a Markov model allows NeSy-MMs to control LLMs

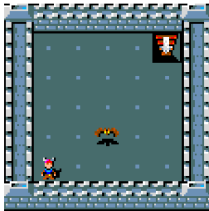


NeSy-MMs approximate safe generative distribution  
while allowing for programming of different constraints

Agents should behave safely  
during exploration and training



Don't go right  
into the lava



Don't go towards  
the monster



Get a key  
before opening  
a locked door

# Modelling policies as Nesy-MMS makes safe RL agents possible

## Previously in our group

Program safety specification in probabilistic logic  
and use conditioned policy for learning and inference

# Modelling policies as Nesy-MMS makes safe RL agents possible

## Previously in our group

Program safety specification in probabilistic logic  
and use conditioned policy for learning and inference

## Limitation

Exact inference is too expensive for large domains  
and only looks ahead one step

# Modelling policies as Nesy-MMS makes safe RL agents possible

## Previously in our group

Program safety specification in probabilistic logic  
and use conditioned policy for learning and inference

## Limitation

Exact inference is too expensive for large domains  
and only looks ahead one step

## Now

NeSy-MMs encode the conditioned policies  
allowing safe policies in large temporal domains



# Neural nets need relational probabilistic reasoning to achieve the guarantees they lack

Relational NeSy-MMs...

...are scalable and differentiable  
neural + logical + probabilistic models

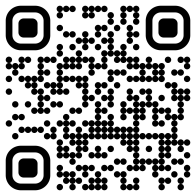
...are logically consistent and intervenable  
while generalising better than other sequential models

...show promise for larger applications  
such as controlled language generation and safe RL



Artificial  
Intelligence  
Viaanderen/Flanders

WASP | WALLENBERG AI  
AUTONOMOUS SYSTEMS  
AND SOFTWARE PROGRAM



Visit personal page



Read our paper