

# Keeping the promises and overcoming the challenges of neurosymbolic AI

Lennert De Smet

How can a robot **safely** bring your favourite pizza  
to your doorstep?



Neural networks

Backpropagation

Symbolic methods

Logic

Neurosymbolic AI (NeSy)



Probability theory

Neural networks

Backpropagation

Symbolic methods

Logic

Probabilistic NeSy



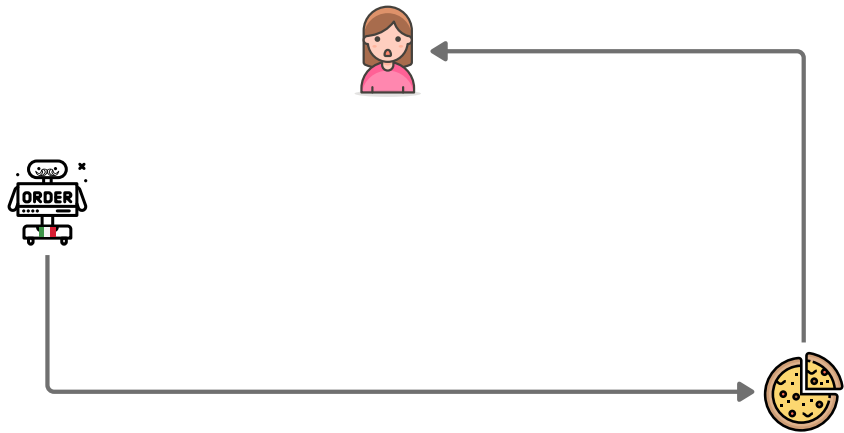
Probability theory



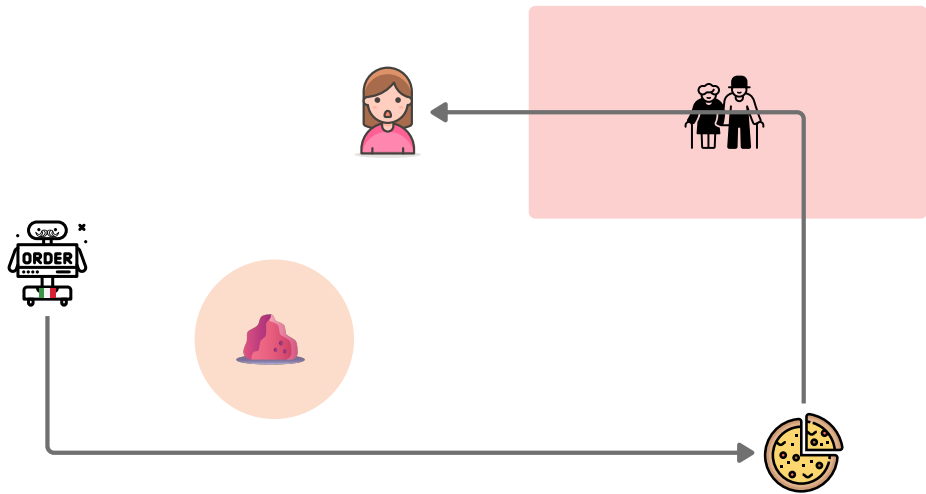














$$\sqrt{(X - R)^2} > 2$$



$$(|X_1 - G_1| > 5) \wedge (|X_2 - G_2| > 3)$$





$$\sqrt{(X - R)^2} > 2$$



$$(|X_1 - G_1| > 5) \wedge (|X_2 - G_2| > 3)$$



- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 How can we still make neurosymbolic AI work?

## 1 What is neurosymbolic AI and what does it promise?

- 1 Neural + Probabilistic + Symbolic AI
- 2 Deep Bayesian networks and neural probabilistic logic
- 3 Guaranteed consistency and improved generalisation

## 2 Where does neurosymbolic AI struggle?

## 3 How can we still make neurosymbolic AI work?

## Neural Models

+ Incredible function approximates

On the left sidewalk there is a young child and on the right sidewalk is an old couple. You are driving on the road and have to take an action, what do you do?



This is a moral and ethical dilemma often presented in discussions about autonomous vehicles and decision-making in critical situations. The question essentially asks you to choose between harming a young child or an old couple, which is a deeply difficult and subjective decision.



## Neural Models

- + Incredible function approximates
- Black box, not robust
- Struggle with reasoning, consistency

On the left sidewalk there is a young child and on the right sidewalk is an old couple. You are driving on the road and have to take an action, what do you do?



This is a moral and ethical dilemma often presented in discussions about autonomous vehicles and decision-making in critical situations. The question essentially asks you to choose between harming a young child or an old couple, which is a deeply difficult and subjective decision.



## Neural Models

- + Incredible function approximates
- Black box, not robust
- Struggle with reasoning, consistency

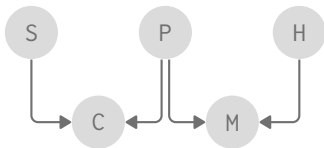
On the left sidewalk there is a young child and on the right sidewalk is an old couple. You are driving on the road and have to take an action, what do you do?



This is a moral and ethical dilemma often presented in discussions about autonomous vehicles and decision-making in critical situations. The question essentially asks you to choose between harming a young child or an old couple, which is a deeply difficult and subjective decision.

## Probabilistic Models

- + Flexible and robust



## Neural Models

- + Incredible function approximates
- Black box, not robust
- Struggle with reasoning, consistency

On the left sidewalk there is a young child and on the right sidewalk is an old couple. You are driving on the road and have to take an action, what do you do?



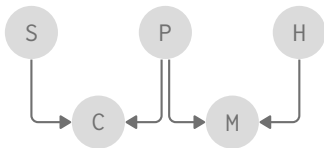
This is a moral and ethical dilemma often presented in discussions about autonomous vehicles and decision-making in critical situations. The question essentially asks you to choose between harming a young child or an old couple, which is a deeply difficult and subjective decision.

## Probabilistic Models

- + Flexible and robust
- Expressivity versus tractability
- No logical reasoning

## Neural Models

- + Incredible function approximates
- Black box, not robust
- Struggle with reasoning, consistency



On the left sidewalk there is a young child and on the right sidewalk is an old couple. You are driving on the road and have to take an action, what do you do?



This is a moral and ethical dilemma often presented in discussions about autonomous vehicles and decision-making in critical situations. The question essentially asks you to choose between harming a young child or an old couple, which is a deeply difficult and subjective decision.

## Deep Probabilistic Models

- + Incredible function approximates
- + Flexible and robust
- Struggle with reasoning, consistency

## Symbolic Models

- + Consistency of logic

## Deep Probabilistic Models

- + Incredible function approximates
- + Flexible and robust
- Struggle with reasoning, consistency

## Symbolic Models

- + Consistency of logic
- Not designed for uncertainty
- Can not deal with “raw” data

## Deep Probabilistic Models

- + Incredible function approximates
- + Flexible and robust
- Struggle with reasoning, consistency

## Probabilistic Neurosymbolic Models

- + Incredible function approximates
- + Flexible and robust
- + Consistency of logical reasoning



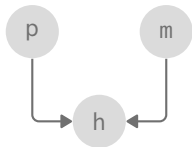
Bayesian networks and probabilistic logic programs  
encode probability distributions

Bayesian networks

Probabilistic logic programs

# Bayesian networks and probabilistic logic programs encode probability distributions

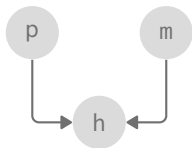
## Bayesian networks



## Probabilistic logic programs

# Bayesian networks and probabilistic logic programs encode probability distributions

## Bayesian networks

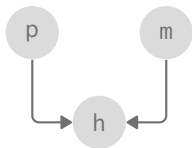


## Probabilistic logic programs

## Defining CPTs

# Bayesian networks and probabilistic logic programs encode probability distributions

## Bayesian networks



## Probabilistic logic programs

```
0.9 :: player((1, 2)).  
0.5 :: monster((1, 2)).
```

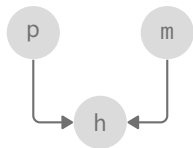
```
hit ← player(L) ∧ monster(L).
```

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

## Defining CPTs

# Bayesian networks and probabilistic logic programs encode probability distributions

## Bayesian networks



## Defining CPTs

## Probabilistic logic programs

```
0.9 :: player((1, 2)).  
0.5 :: monster((1, 2)).
```

```
hit ← player(L) ∧ monster(L).
```

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

## Writing logic programs

Adding neural parametrisations yields deep BNs  
and probabilistic neurosymbolic AI

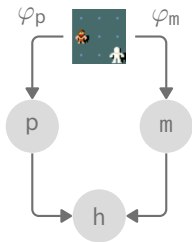
Deep Bayesian networks

Neural Probabilistic logic programs

# Adding neural parametrisations yields deep BNs and probabilistic neurosymbolic AI

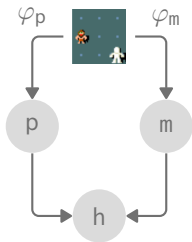
Deep Bayesian networks

Neural Probabilistic logic programs



# Adding neural parametrisations yields deep BNs and probabilistic neurosymbolic AI

## Deep Bayesian networks



## Neural Probabilistic logic programs

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

$\text{hit}(\text{img}) \leftarrow \text{player}(\text{img}, L) \wedge \text{monster}(\text{img}, L).$

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



# Impose consistency on neural models by conditioning on logical statements

Take a neural distribution  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

# Impose consistency on neural models by conditioning on logical statements

Take a neural distribution  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

$$p_{\theta}(\omega \mid \varphi) = \frac{p(\varphi \mid \omega) p_{\theta}(\omega)}{p(\varphi)}$$

# Impose consistency on neural models by conditioning on logical statements

Take a neural distribution  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

$$p_{\theta}(\omega \mid \varphi) = \frac{p(\varphi \mid \omega) p_{\theta}(\omega)}{p(\varphi)}$$

Then Bayes' rule give us new predictions  
that are **consistent** with the logic formula

- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 How can we still make neurosymbolic AI work?

Probabilistic inference is **#P**-hard  
and adding logic complicates things further

Take a neural distribution  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

Probabilistic inference is **#P**-hard  
and adding logic complicates things further

Take a neural distribution  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

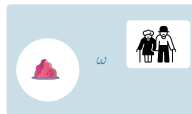
$$p(\varphi) = \int_{\Omega} \mathbb{1}_{\omega \models \varphi} p_{\theta}(\omega) d\omega$$

Probabilistic inference is **#P**-hard  
and adding logic complicates things further

Take a **neural distribution**  $p_{\theta}(\omega)$   
and a logic formula  $\varphi$  over  $\omega$

$$p(\varphi) = \int_{\Omega} \mathbb{1}_{\omega \models \varphi} p_{\theta}(\omega) d\omega$$

Remember example



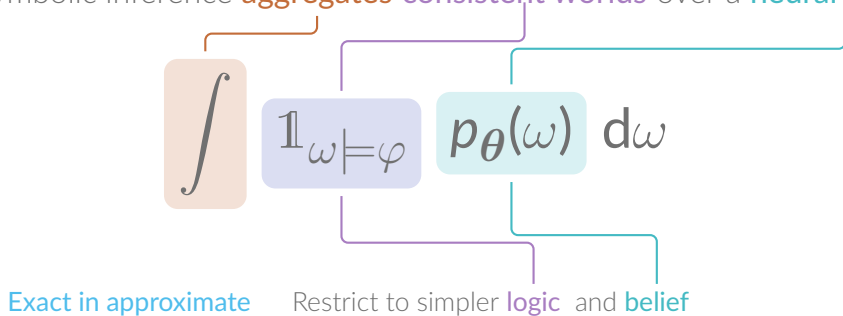
Neurosymbolic inference aggregates consistent worlds over a neural weight

The diagram illustrates the components of the neurosymbolic inference formula  $\int \mathbb{1}_{\omega \models \varphi} p_{\theta}(\omega) d\omega$ . It features three colored boxes: an orange box for the integral symbol  $\int$ , a purple box for the consistency indicator  $\mathbb{1}_{\omega \models \varphi}$ , and a teal box for the probability density  $p_{\theta}(\omega)$ . The differential  $d\omega$  is placed to the right of the teal box. Three brackets connect the text above to the boxes: an orange bracket connects 'aggregates' to the integral symbol, a purple bracket connects 'consistent worlds' to the consistency indicator, and a teal bracket connects 'neural weight' to the probability density.

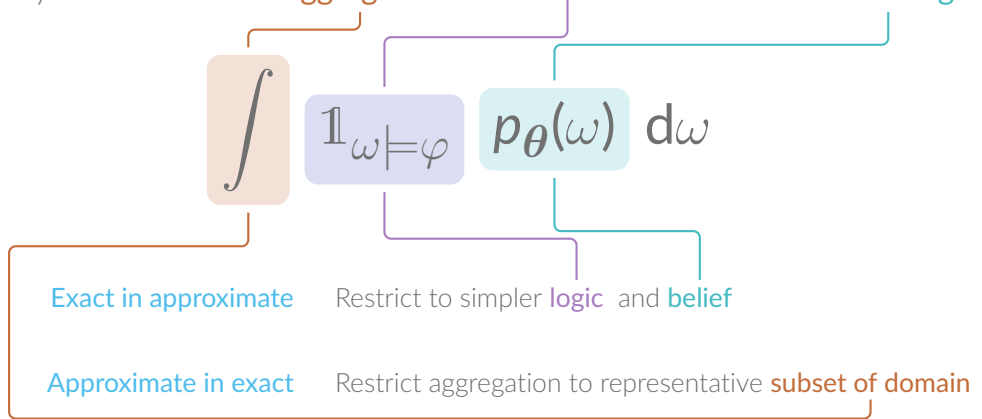
$$\int \mathbb{1}_{\omega \models \varphi} p_{\theta}(\omega) d\omega$$



Neurosymbolic inference aggregates consistent worlds over a neural weight



Neurosymbolic inference **aggregates consistent worlds** over a **neural weight**



- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 How can we still make neurosymbolic AI work?

- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 **How can we still make neurosymbolic AI work?**
  - 1 Neurosymbolic models for sequential data
  - 2 Turning hallucinations into consistency and learning generalising models
  - 3 Constraining language models and safe reinforcement learning

Yes, we can make neurosymbolic AI work  
on challenging sequential data!

**Problem** Transformers are amazing at learning from sequential data  
but often hallucinate and struggle with consistency

Yes, we can make neurosymbolic AI work on challenging sequential data!

**Problem** Transformers are amazing at learning from sequential data but often hallucinate and struggle with consistency

**Challenge** Existing methods for neurosymbolic inference were not able to scale to sequential data

Yes, we can make neurosymbolic AI work on challenging sequential data!

**Problem** Transformers are amazing at learning from sequential data but often hallucinate and struggle with consistency

**Challenge** Existing methods for neurosymbolic inference were not able to scale to sequential data

**Solution** Our neurosymbolic Markov models (NeSy-MMs) combine sequential probabilistic models with symbolic logic



NeSy



Sequential



Relational



Discrete and continuous

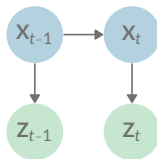


Neural + logical



Discriminative and generative





NeSy

HMM



Sequential



Relational



Discrete and continuous



Neural + logical



Discriminative and generative





NeSy



Sequential



Relational



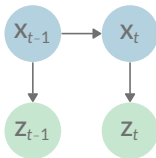
Discrete and continuous



Neural + logical



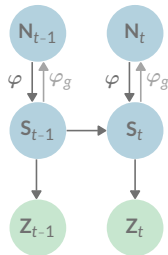
Discriminative and generative



HMM



Our solution

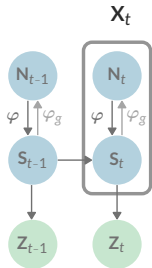


NeSy-MMs



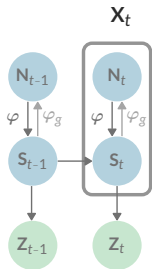
Particle filters scale sequential probabilistic inference  
in discrete-continuous domains

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



Particle filters scale sequential probabilistic inference  
in discrete-continuous domains

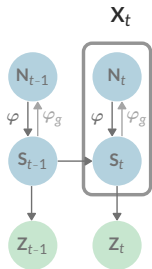
$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



1 Recursively draw samples

Particle filters scale sequential probabilistic inference in discrete-continuous domains

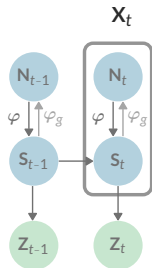
$$p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{Z}_{0:t+1}) \propto \int p_{\varphi}(\mathbf{Z}_{t+1} | \mathbf{X}_{t+1}) \cdot p_{\varphi}(\mathbf{X}_{t+1} | \mathbf{x}_t) \cdot p_{\varphi}(\mathbf{x}_t | \mathbf{Z}_{0:t}) d\mathbf{x}_t$$



- 1 **Recursively** draw samples
- 2 **Transition** recursive samples

# Particle filters scale sequential 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 **Recursively** draw samples
- 2 **Transition** recursive samples
- 3 **Resample** transitioned samples with observation

Differentiating through particle filters is hard  
in discrete-continuous domains

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

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

Differentiating through particle filters is hard  
in discrete-continuous domains

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

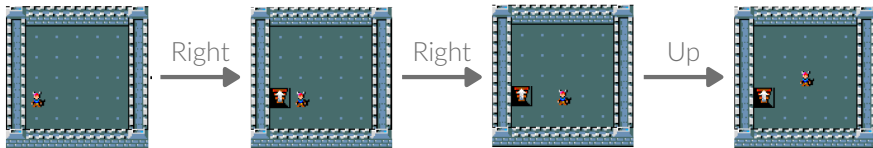
**Problem** Resampling is not a differentiable operation  
as it is the same as sampling from a discrete distribution

**Solution** Existing work tackles continuous resampling gradients  
and we add a solution for discrete resampling gradients

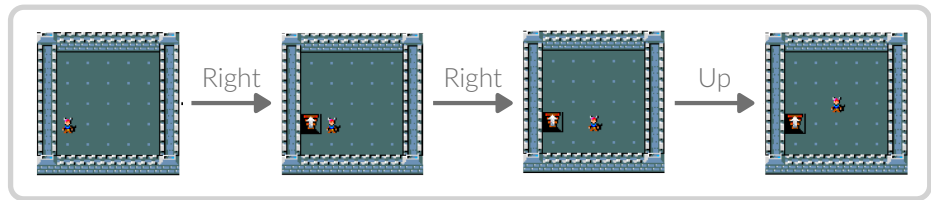


- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 How can we still make neurosymbolic AI work?
  - 1 Neurosymbolic models for sequential data
  - 2 Turning hallucinations into consistency and learning generalising models
  - 3 Constraining language models and safe reinforcement learning

## Generative data

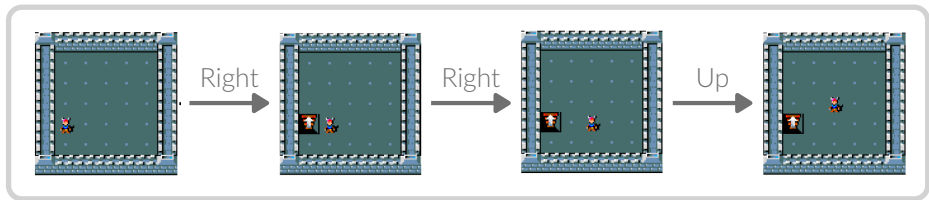


## Generative data

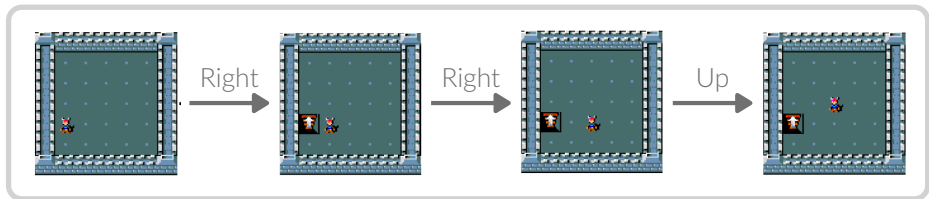


Model

## Generative data



## Generative data



R, D, L, U, L

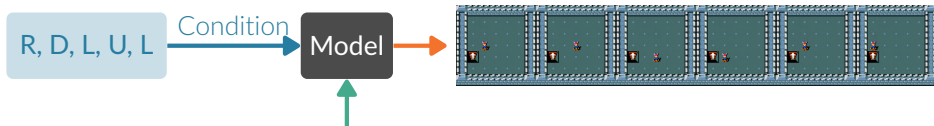
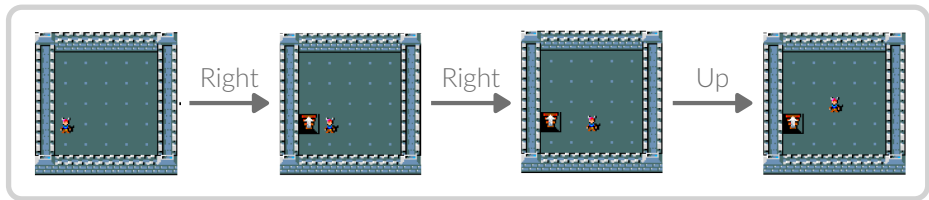
Condition

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

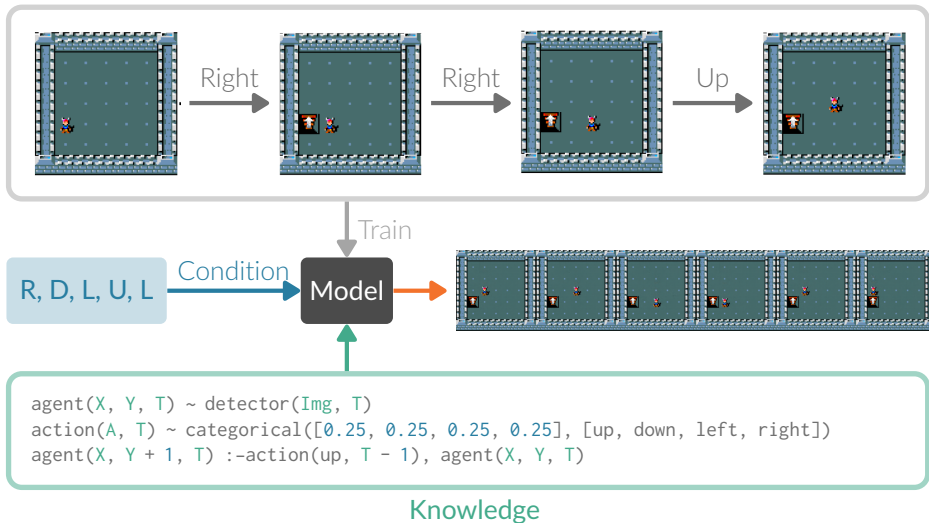
## Generative data



```
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

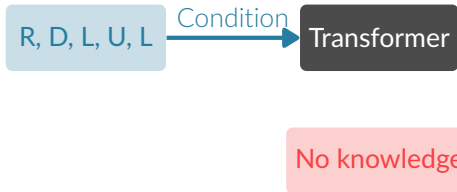
Transformer



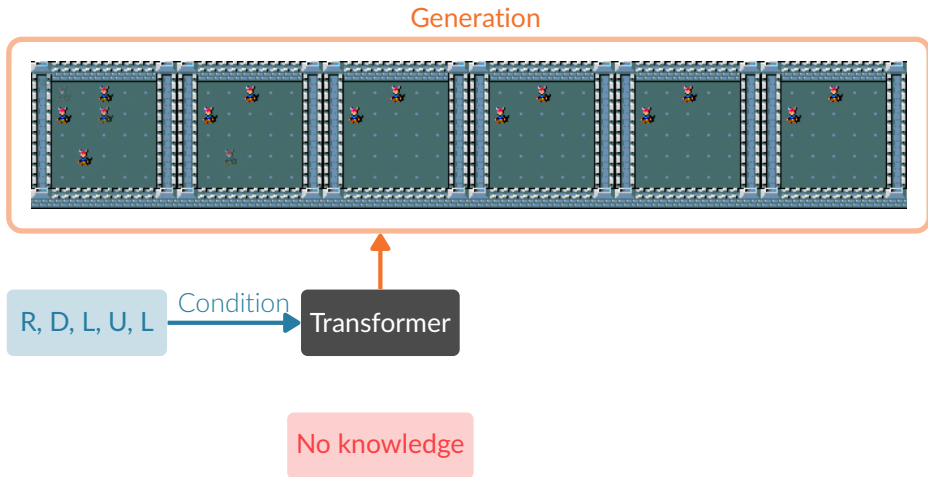
Transformers do not generate logically consistent images



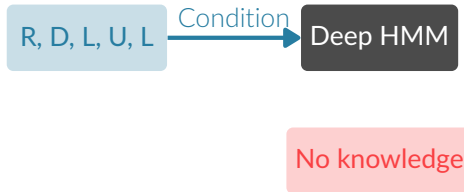
Transformers do not generate logically consistent images



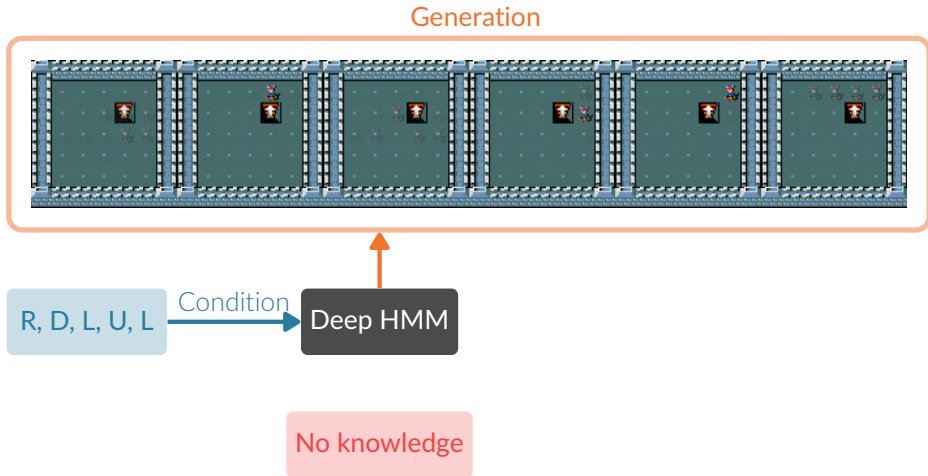
Transformers do not generate logically consistent images



Deep HMMs do not generate logically consistent images



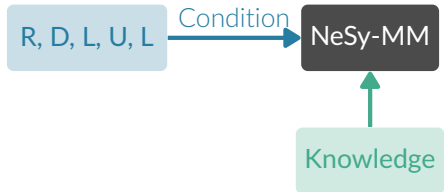
Deep HMMs do not generate logically consistent images



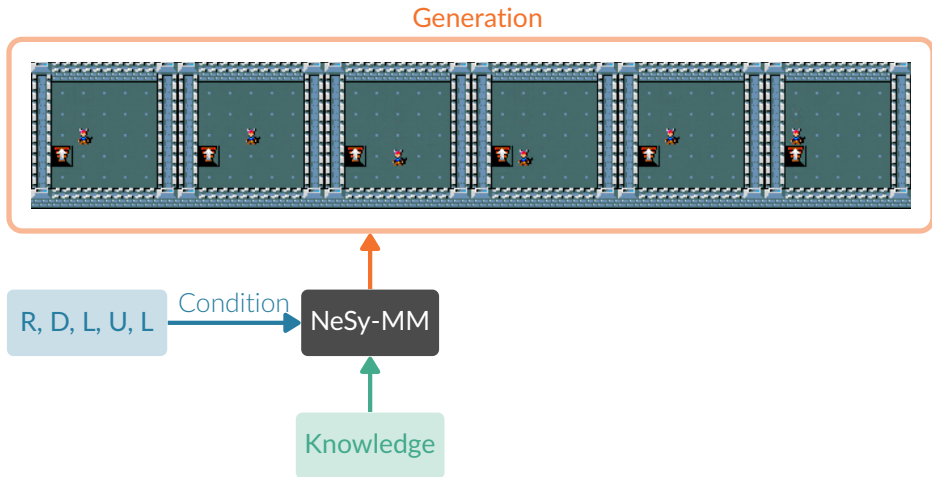
NeSy-MMs do generate logically consistent images



NeSy-MMs do generate logically consistent images

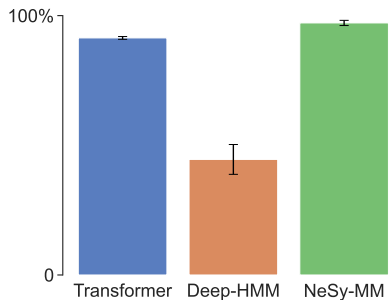


NeSy-MMs do generate logically consistent images



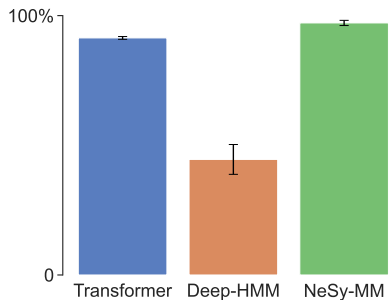


# Quantifying logical consistency with reconstruction accuracy

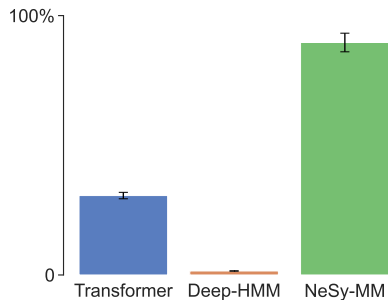


$5 \times 5$

# Quantifying logical consistency with reconstruction accuracy

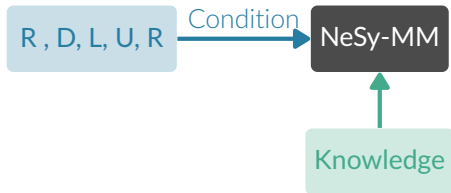


$5 \times 5$

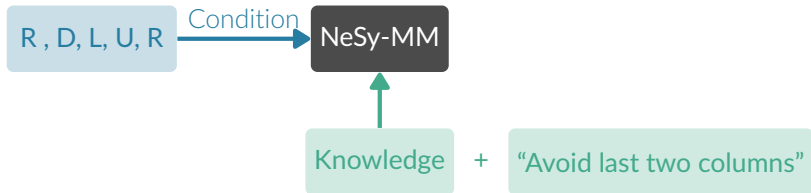


$10 \times 10$

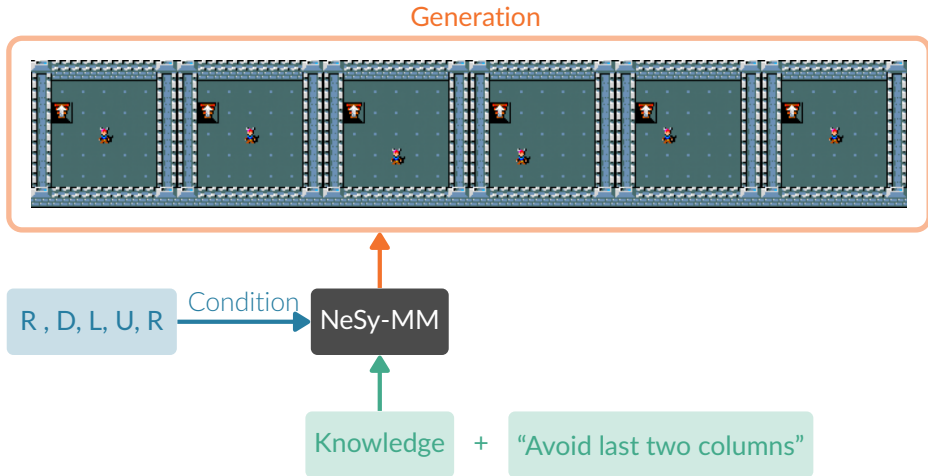
Knowledge can be changed at any time  
without having to retrain any models



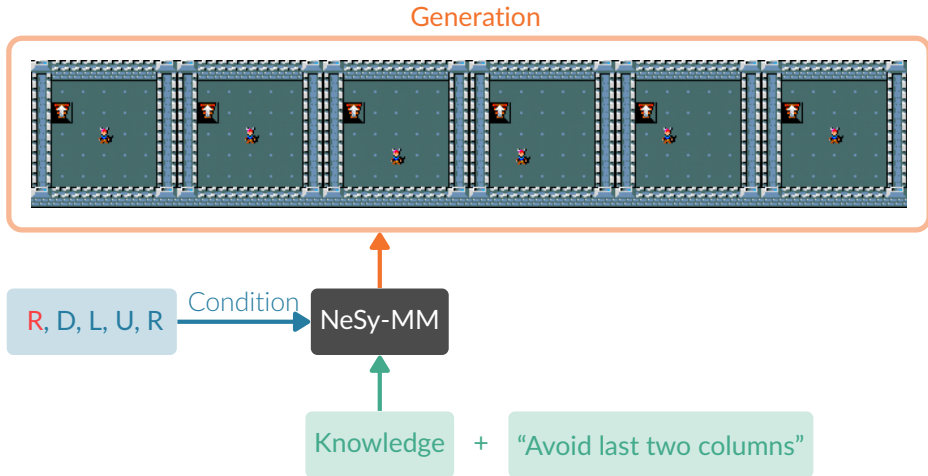
Knowledge can be changed at any time  
without having to retrain any models



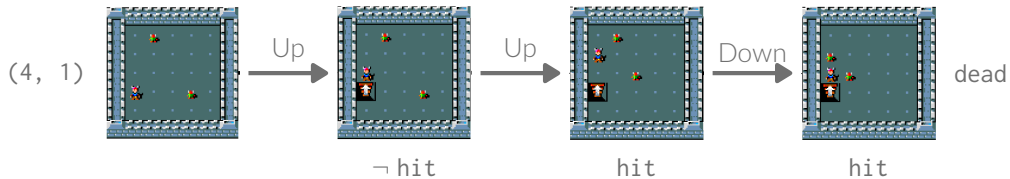
Knowledge can be changed at any time  
without having to retrain any models



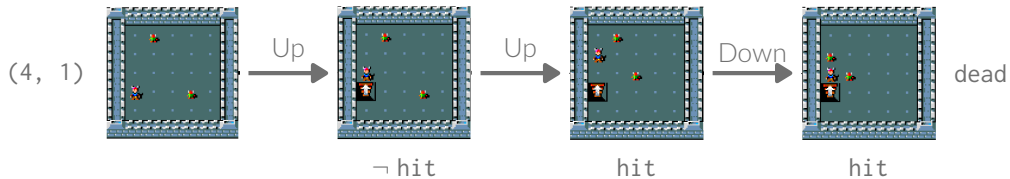
Knowledge can be changed at any time  
without having to retrain any models



## Discriminative symbolic data



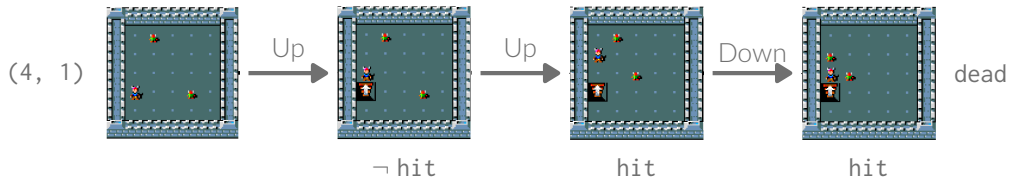
## Discriminative symbolic data



Model

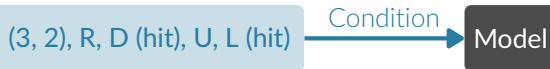
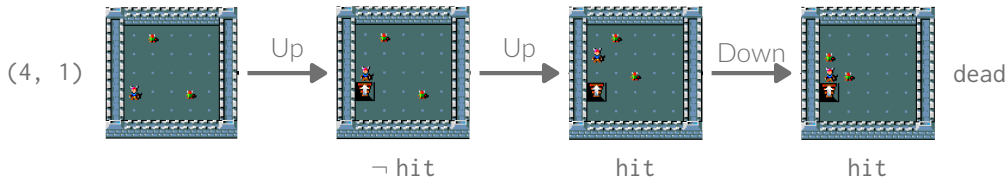


## Discriminative symbolic data



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

## Discriminative symbolic data

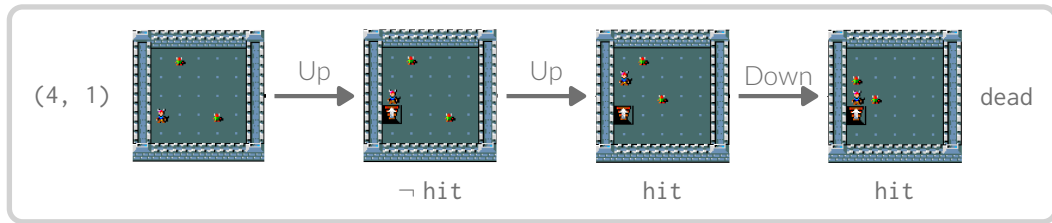


```

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

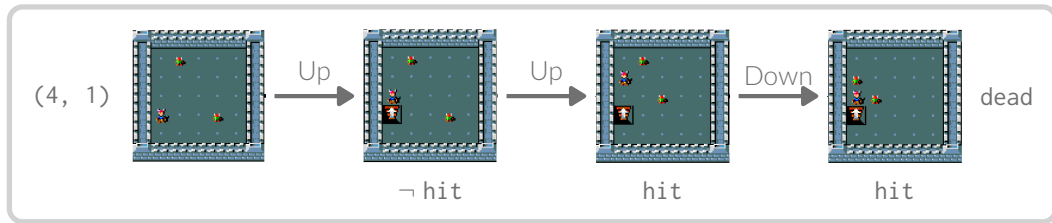


```

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

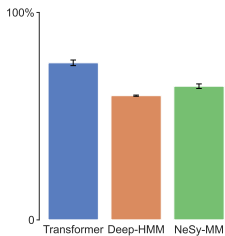


```

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?

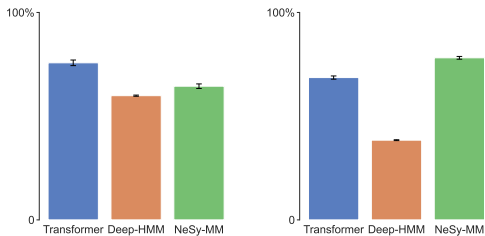


$10 \times 10$

10 steps

1 enemy

# Which model learns a generalisable representation?



$10 \times 10$

10 steps

1 enemy

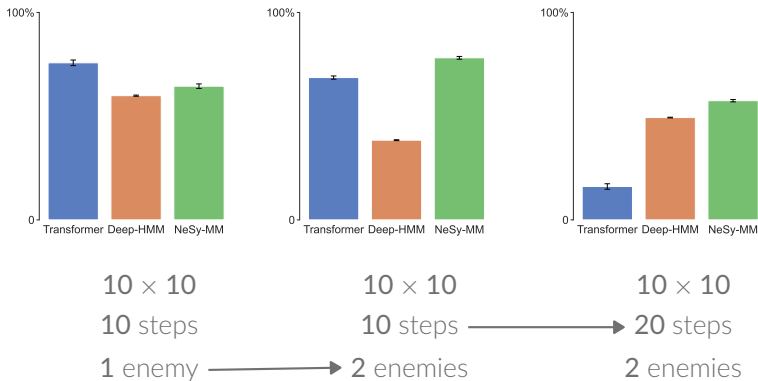
$10 \times 10$

10 steps

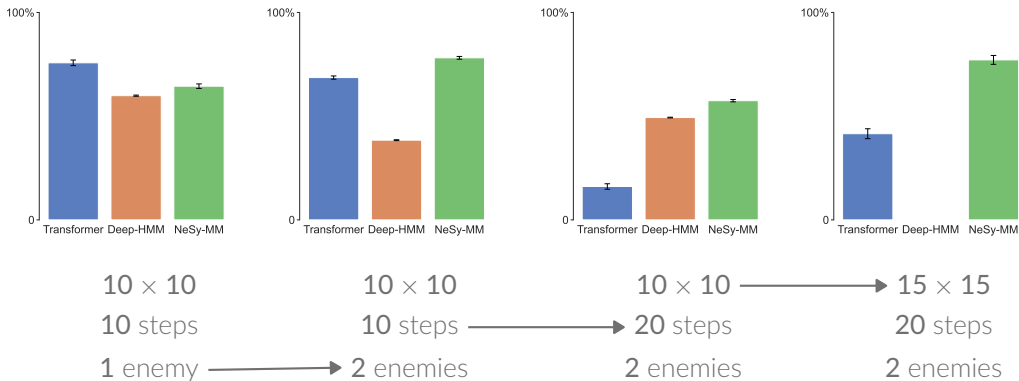
2 enemies



# Which model learns a generalisable representation?

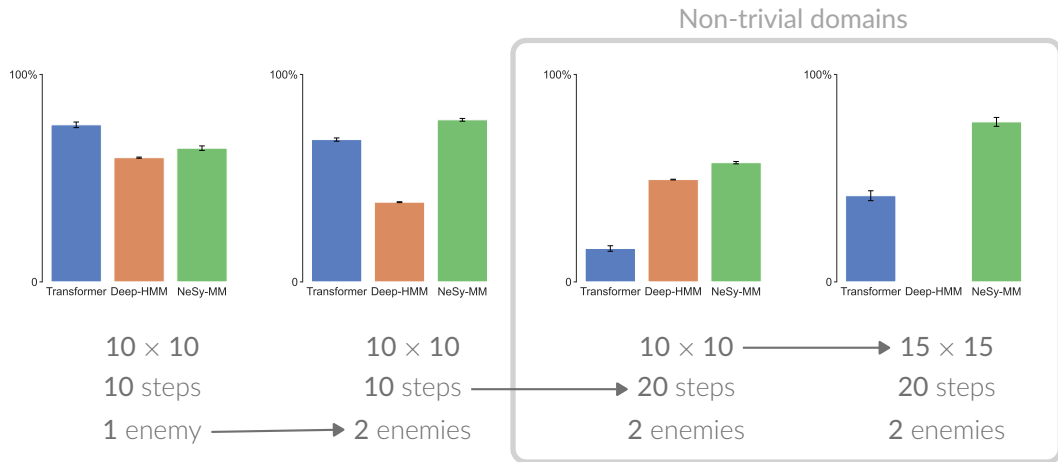


# Which model learns a generalisable representation?





# Which model learns a generalisable representation?



- 1 What is neurosymbolic AI and what does it promise?
- 2 Where does neurosymbolic AI struggle?
- 3 **How can we still make neurosymbolic AI work?**
  - 1 Neurosymbolic models for sequential data
  - 2 Turning hallucinations into consistency and learning generalising models
  - 3 **Constraining language models and safe reinforcement learning**

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?

Can we change constraints without retraining?

# Safe language generation is important but hard guarantee

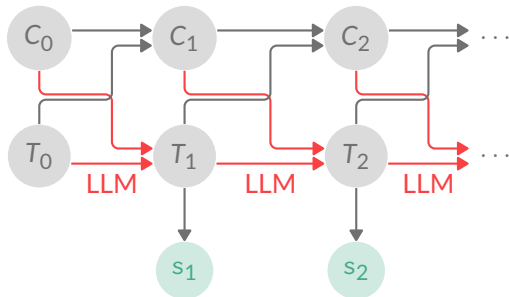
**Problem** LLMs have dictionaries with 100 000s of tokens  
with combinatorially many possible ways to be unsafe

**Questions** Is training for safety enough?

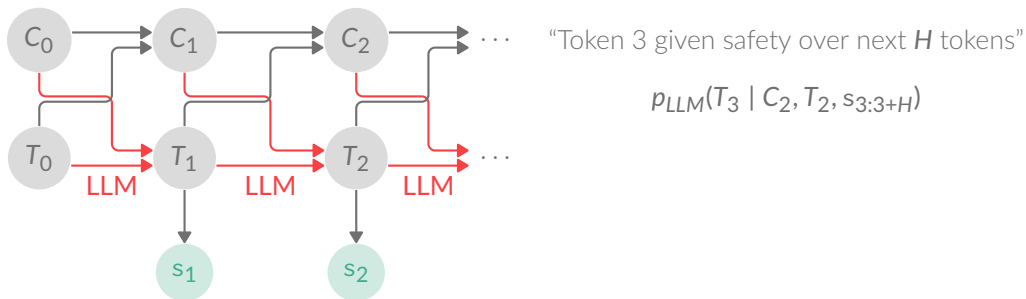
Can we change constraints without retraining?

Does changing predictions influence performance?

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

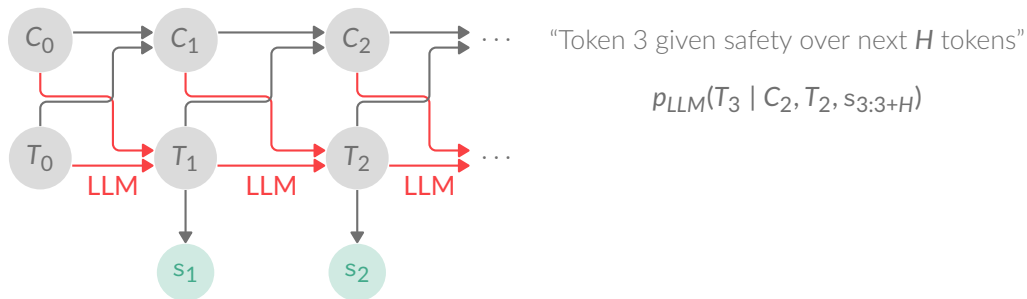


Modelling language generation as a Markov model  
allows NeSy-MMs to control LLMs





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

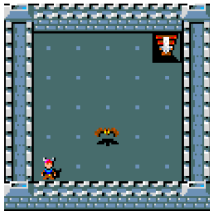


NeSy-MMs approximate safe generative distribution  
while allowing for programming of different constraints

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

Neurosymbolic AI is one avenue to consistent AI  
but there is still much work to do

**Potential**

Neurosymbolic AI models can learn from data  
and guarantee logical consistency

Neurosymbolic AI is one avenue to consistent AI  
but there is still much work to do

**Potential**

Neurosymbolic AI models can learn from data  
and guarantee logical consistency

**Scalability**

Neurosymbolic AI models struggle with large domains  
and require a lot of data to learn

Neurosymbolic AI is one avenue to consistent AI  
but there is still much work to do

**Potential**

Neurosymbolic AI models can learn from data  
and guarantee logical consistency

**Scalability**

Neurosymbolic AI models struggle with large domains  
and require a lot of data to learn

**Solutions exist**

Large sequential domains become viable  
by exploiting approximate inference techniques



Neurosymbolic AI models will deliver your pizza  
hot and safely to your door



Visit personal page



Read our paper