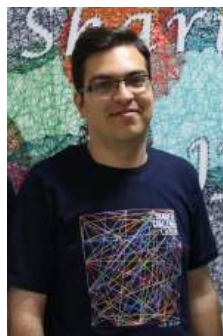


# Semantically-informed Hierarchical Event Modeling



**Shubhashis Roy**  
**Dipta**



**Mehdi Rezaee**



**Frank Ferraro**



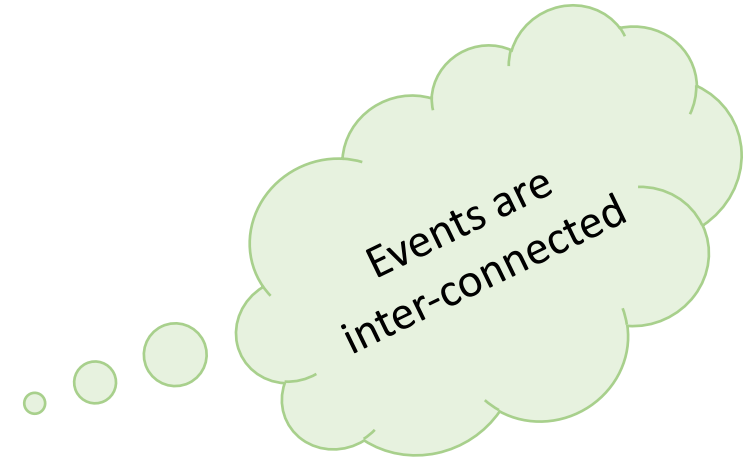
Supported  
by:



Complex events and situations can  
be hierarchical.

This hierarchy presents difficulties.

# Difficulty of Complex Events

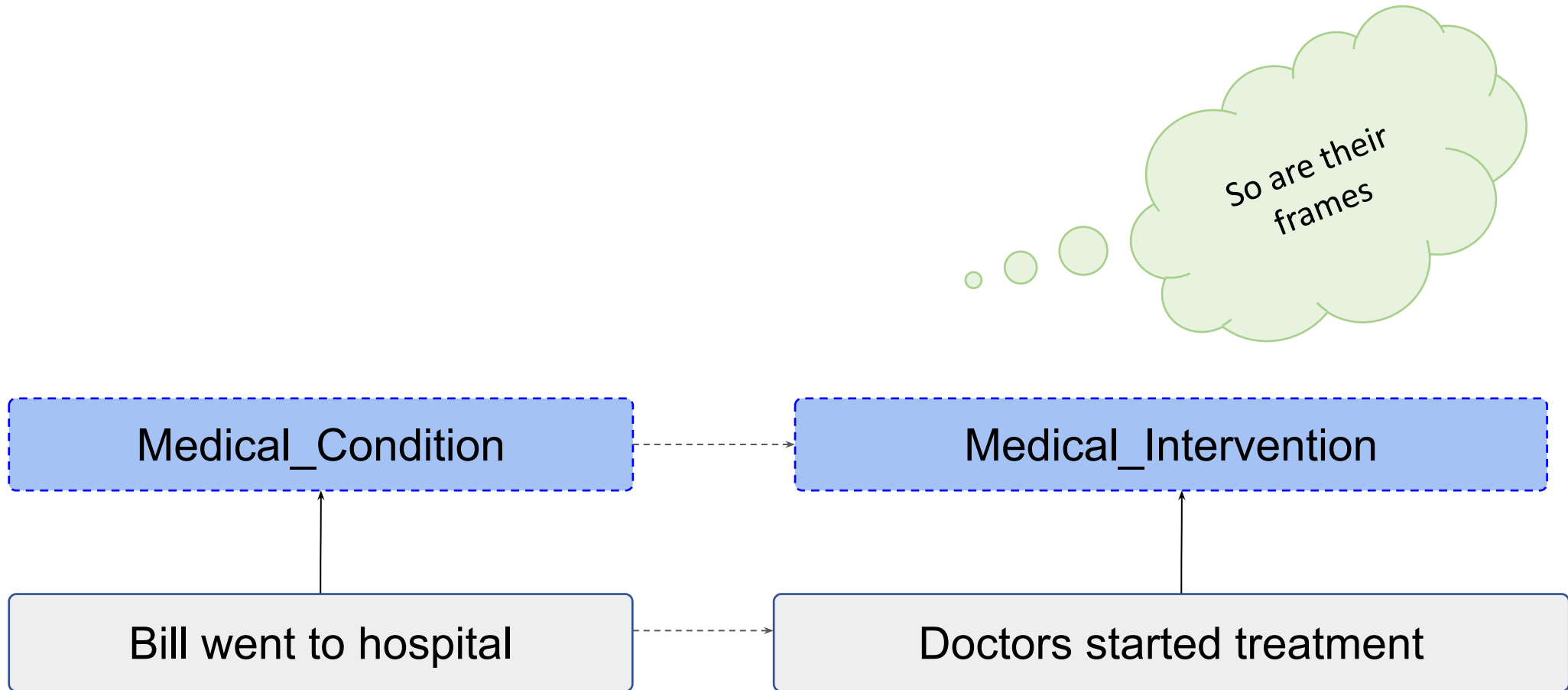


Bill went to hospital

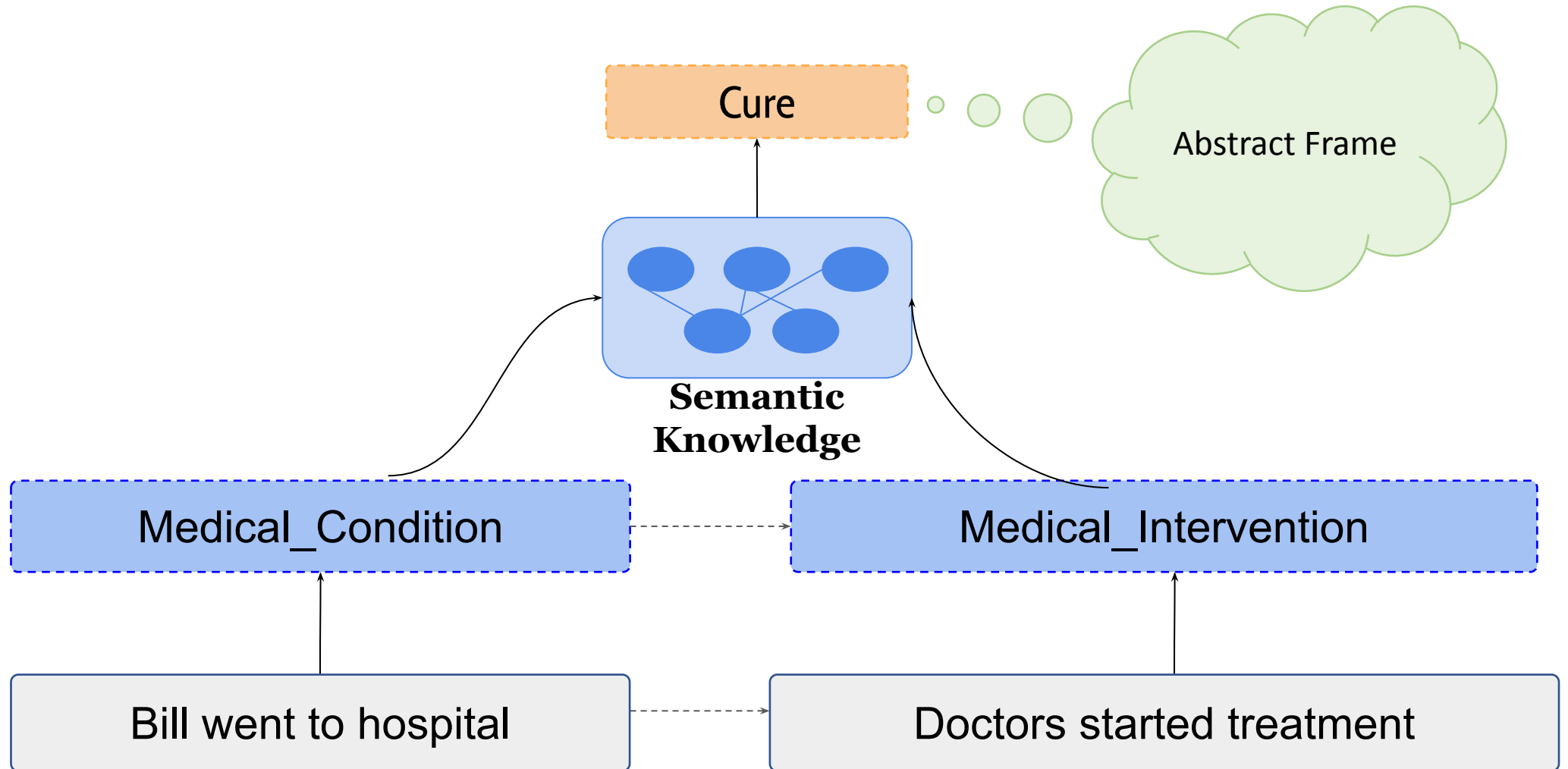


Doctors started treatment

# Difficulty of Complex Events



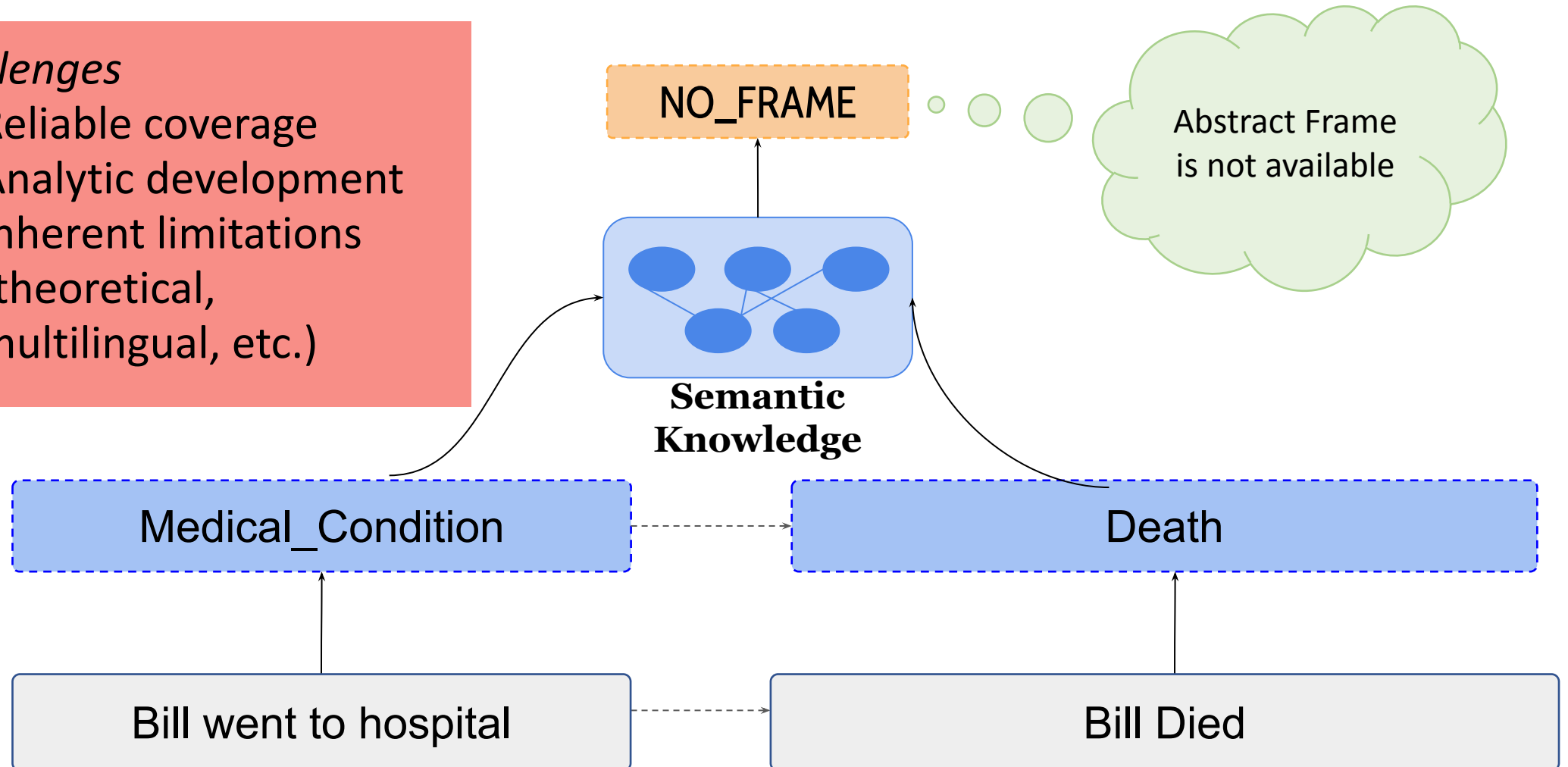
# Semantic Associations can help...



# ...unless you don't have them...

## Challenges

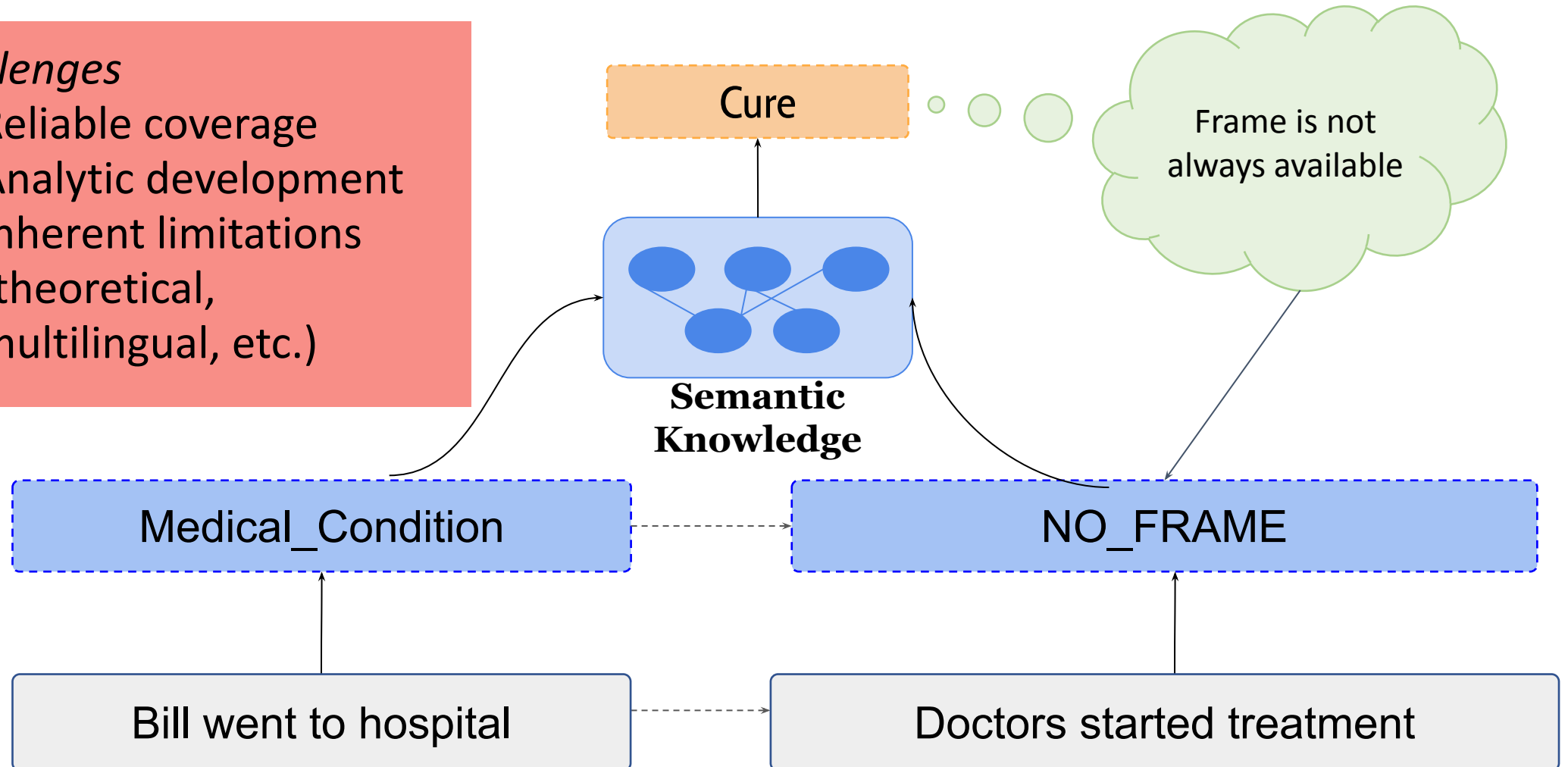
- Reliable coverage
- Analytic development
- Inherent limitations (theoretical, multilingual, etc.)



# ...or if you are missing “lower” semantics

## Challenges

- Reliable coverage
- Analytic development
- Inherent limitations (theoretical, multilingual, etc.)



Complex events and situations are hierarchical. How can we better capture this hierarchy...

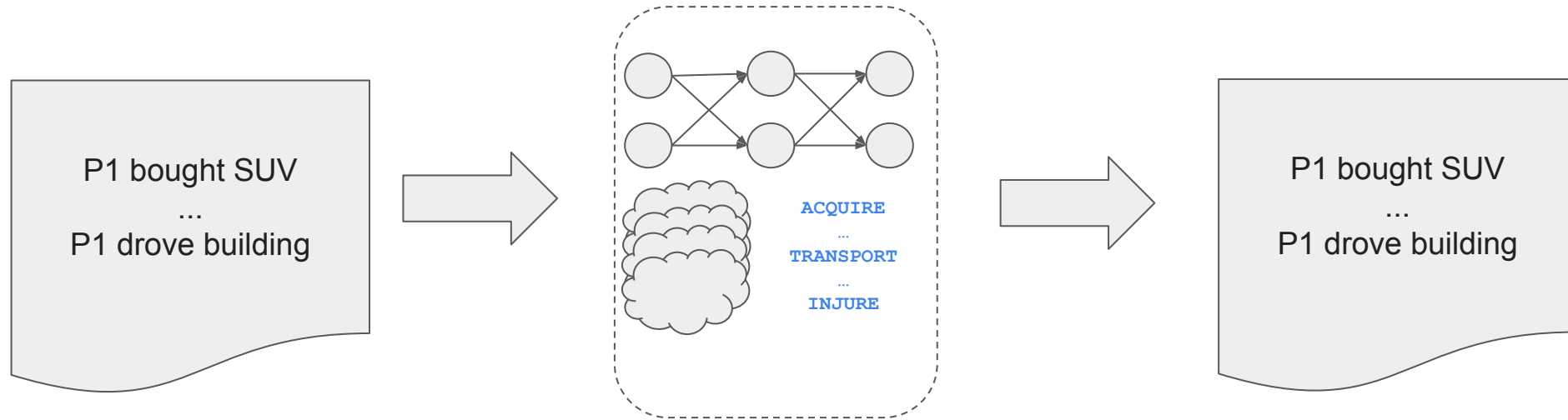
- via modeling improvements?
- utilizing (existing) semantic resources?



# We introduce SHEM:

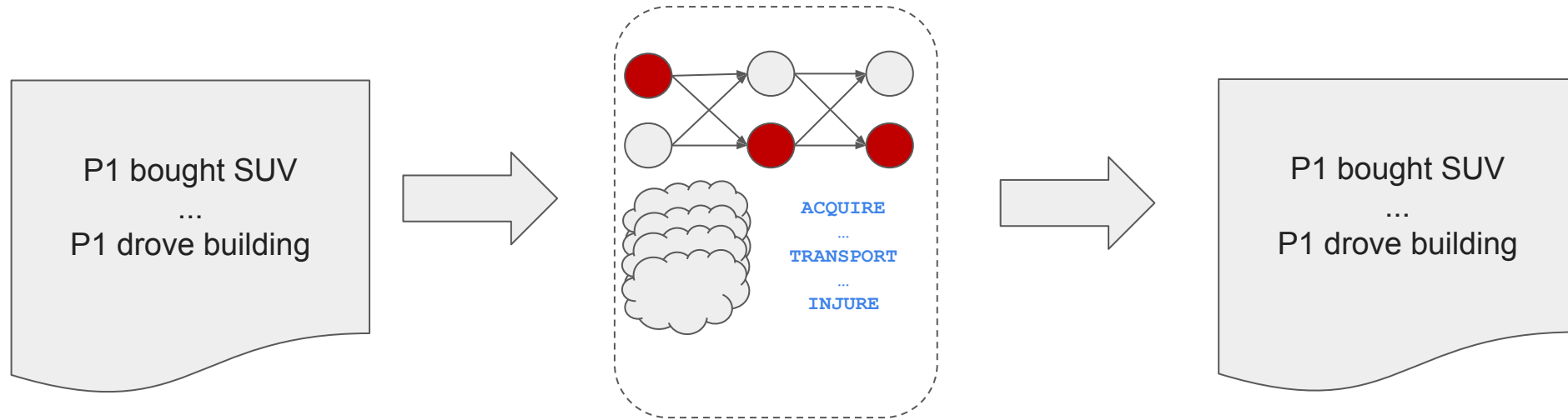
- A Semi-supervised, Hierarchical Event learning Model
- Use of existing FrameNet resource for extracting side knowledge and abstract concept about the events
- Hierarchical model with combination of InfoNCE loss to provide better event representation
- Our model shows better performance in multiple tasks
  - What event comes next
  - Generating missing events
  - Identify similar or related events

# Encoding into Primitives and Back-Again



A neural generative model that uses partially-observed, **semantic (ontology-based) knowledge** to explain and predict events

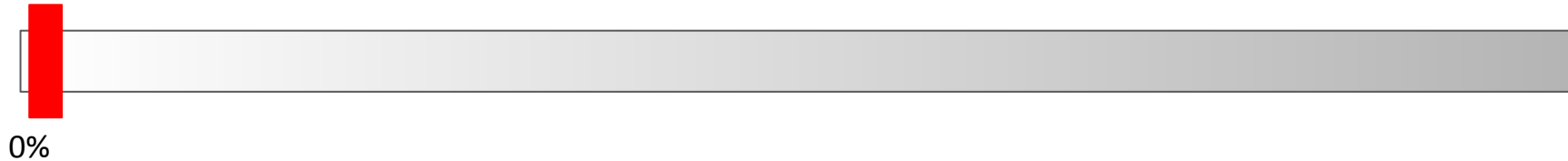
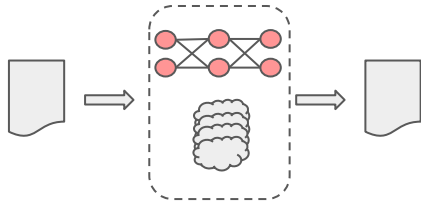
# Encoding into Primitives and Back-Again



A neural generative model that uses partially-observed, **semantic (ontology-based) knowledge** to explain and predict events

Represent **latent variables** with Gumbel-Softmax and ***softly*** inject the information into them

# Encoding into Primitives and Back-Again

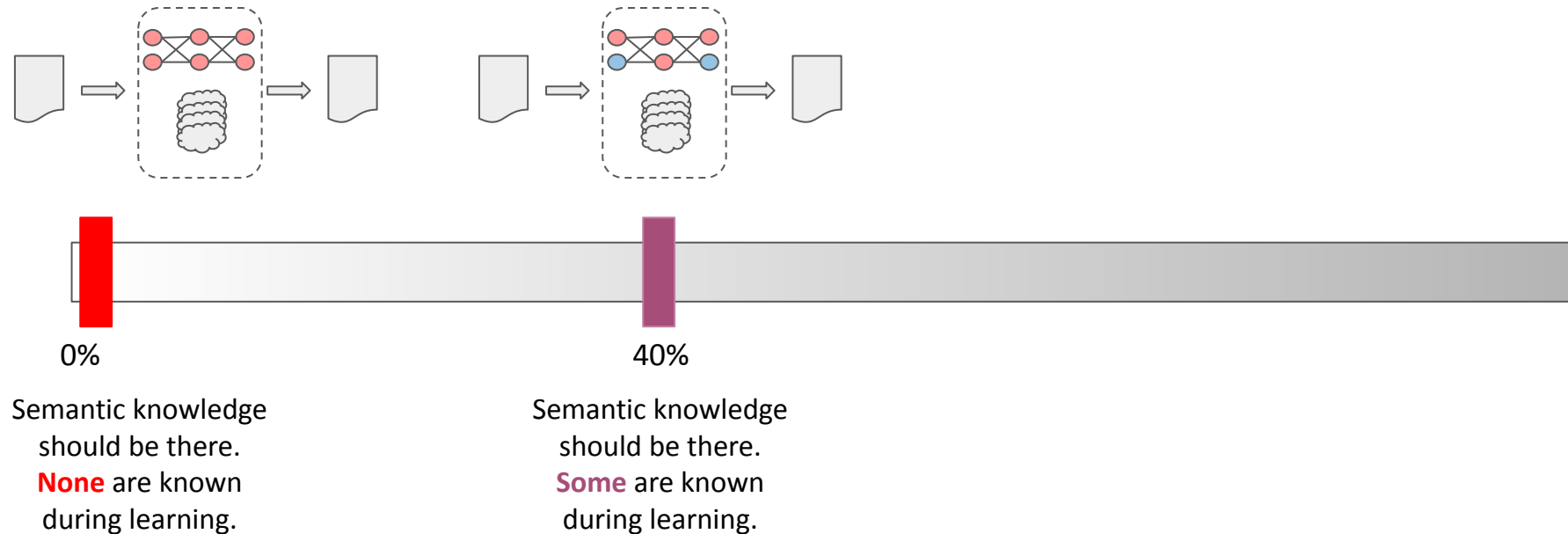


Semantic knowledge  
should be there.  
**None** are known  
during learning.

A neural generative model that uses  
partially-observed, **semantic  
(ontology-based) knowledge**  
to explain and predict events

Represent **latent variables**  
with Gumbel-Softmax and  
**softly** inject the information  
into them

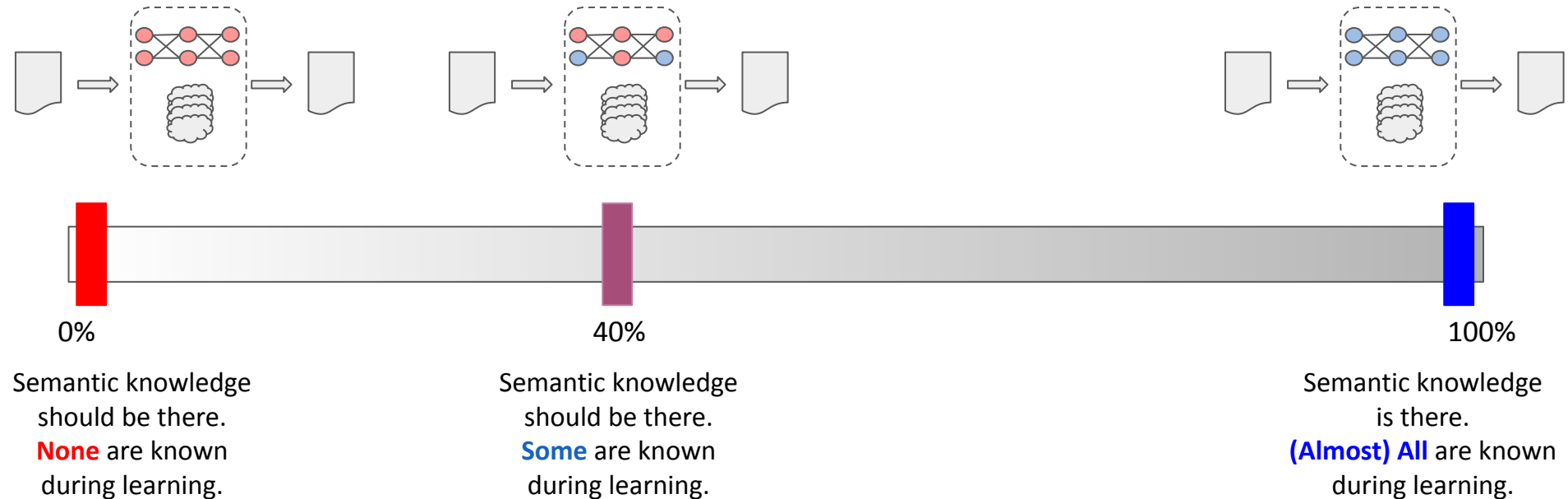
# Encoding into Primitives and Back-Again



A neural generative model that uses partially-observed, **semantic (ontology-based) knowledge** to explain and predict events

Represent **latent variables** with Gumbel-Softmax and **softly** inject the information into them

# Encoding into Primitives and Back-Again

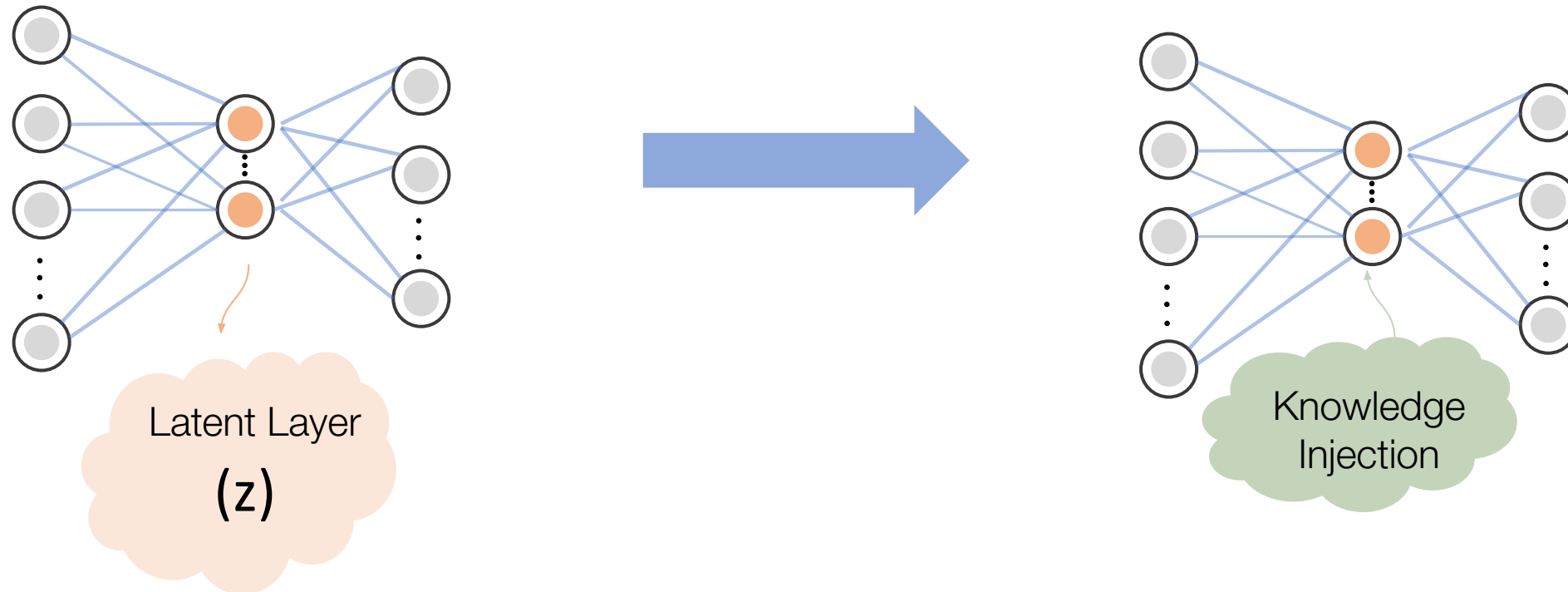


A neural generative model that uses partially-observed, **semantic (ontology-based) knowledge** to explain and predict events

Represent **latent variables** with Gumbel-Softmax and **softly** inject the information into them

# Previous Solutions

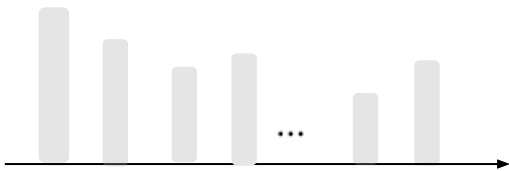
- HAQAE (Weber et al., 2018): latent tree-based model
- SSDVAE (Rezaee and Ferraro, 2021): latent variable method, injecting semantic information to the discrete latent layer parameters.



# Previous Solutions

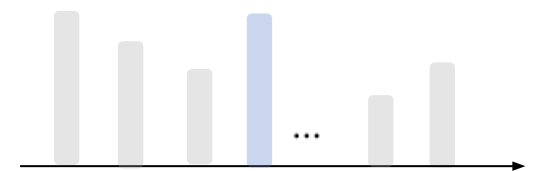
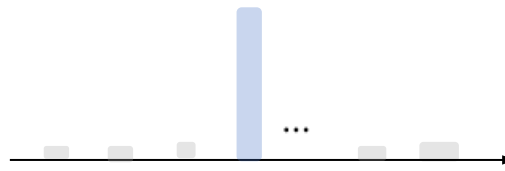
- HAQAE (Weber et al., 2018): latent tree-based model
- SSDVAE (Rezaee and Ferraro, 2021): latent variable method, injecting semantic information to the discrete latent layer parameters.

$\pi$ : network-computed logits



+

$f$ : (scaled) indicator of external knowledge

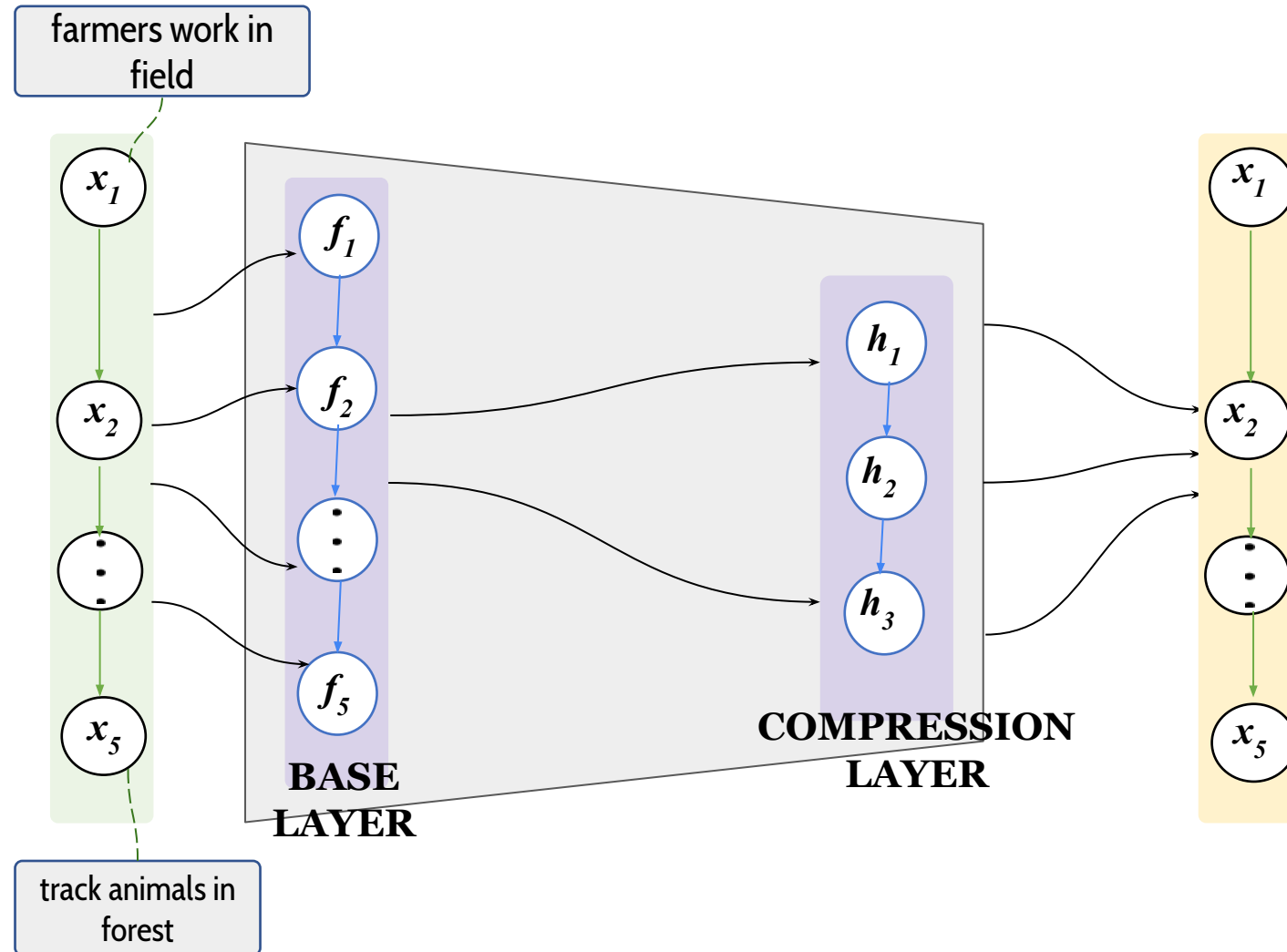




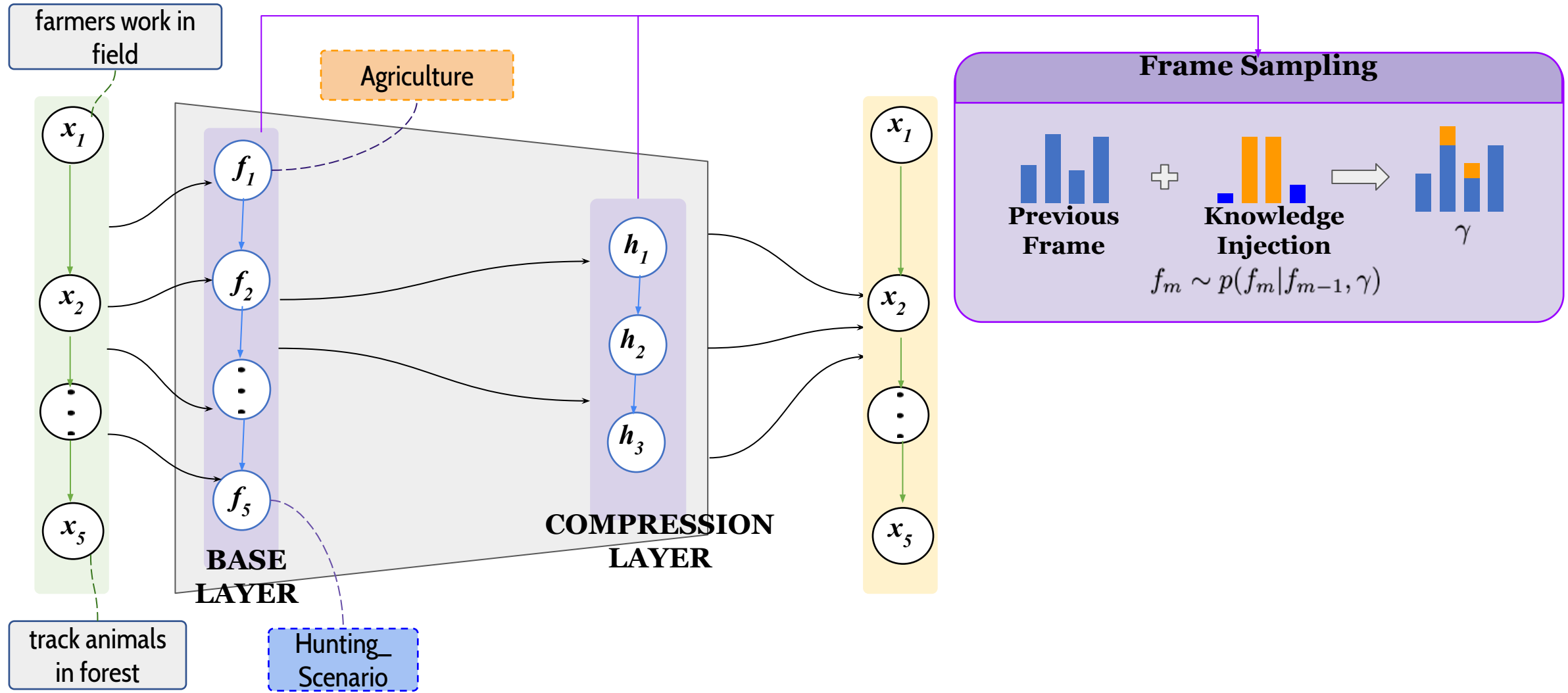
# Our Model Description

- Two Layers
  - Base Layer
    - Token Encoder
    - Latent Variable Classification
    - Token Decoder
  - Compression Layer
    - Semantic Knowledge Extraction
    - Sequence Encoder
    - Latent Variable Compression
    - Token Decoder

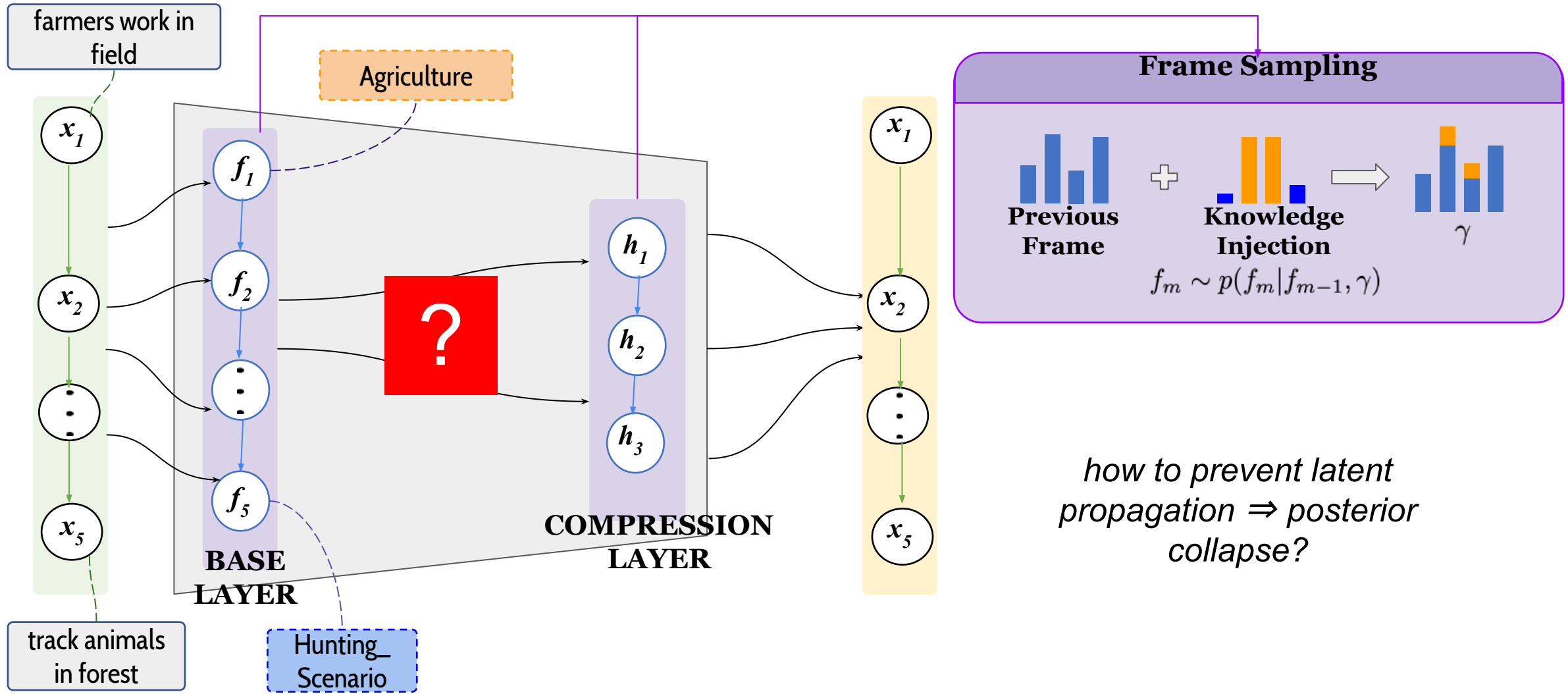
# Capturing Joint Hierarchy: Multi-layer Encoder-Decoder



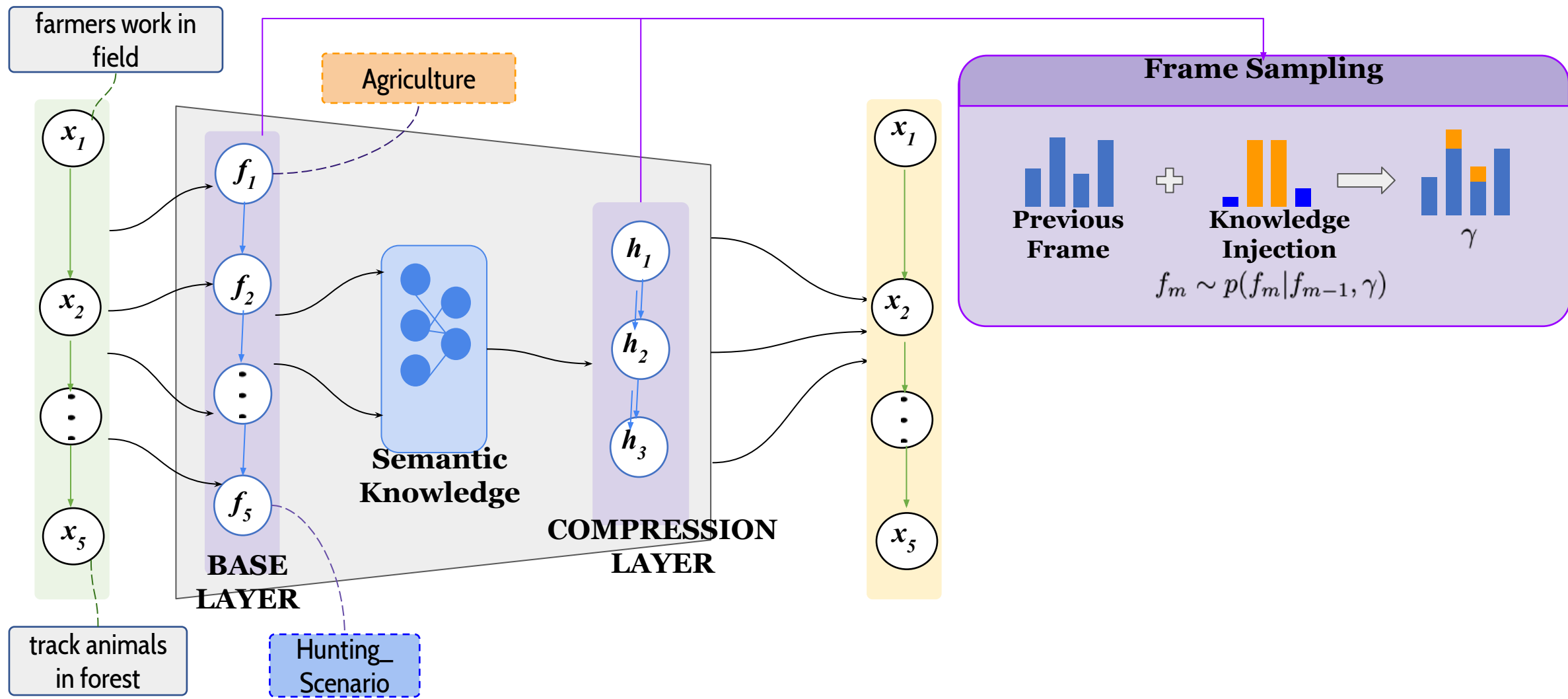
# Capturing Joint Hierarchy: Inject Partially Observable Knowledge



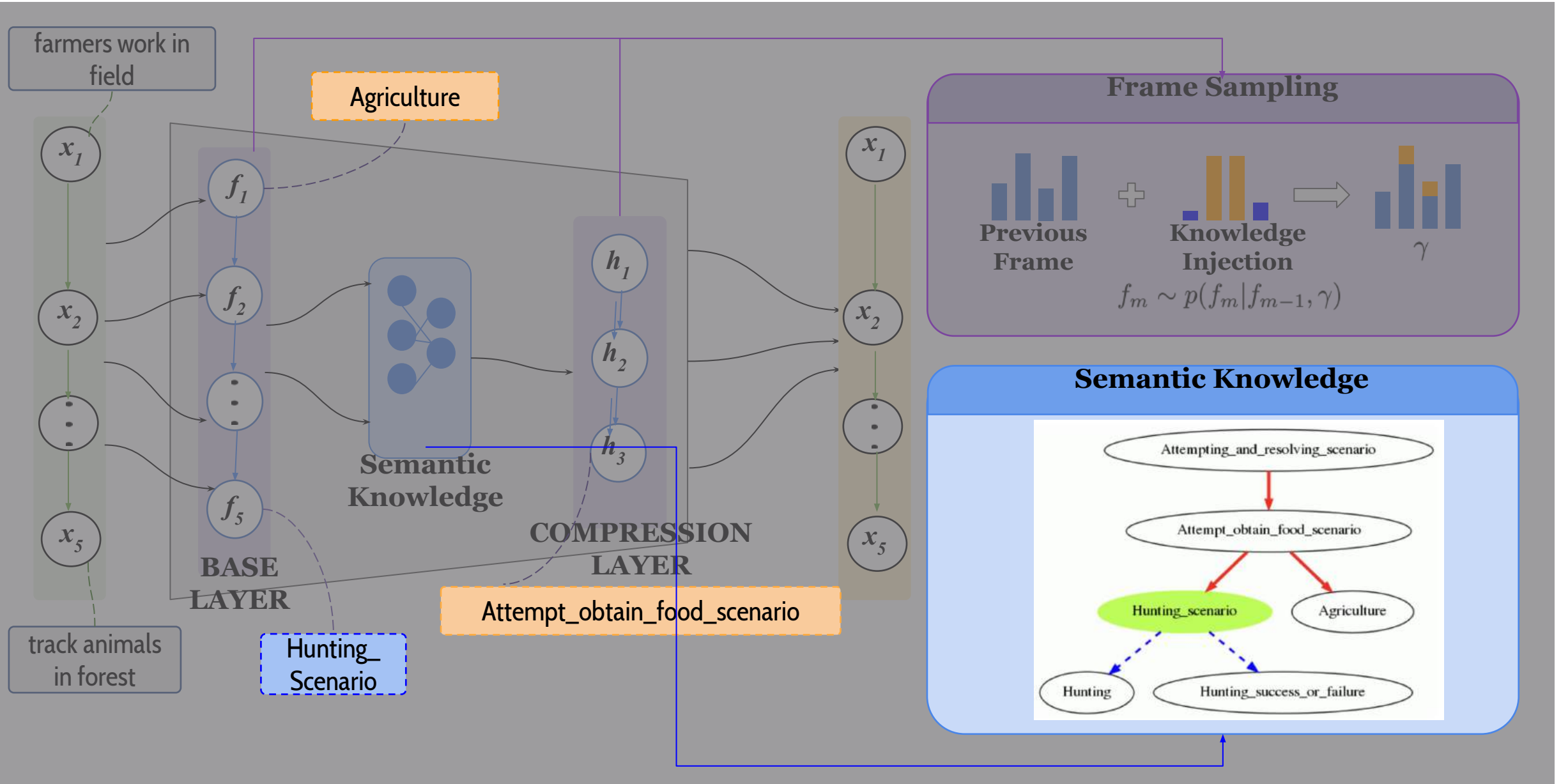
# Capturing Joint Hierarchy: Inject Partially Observable Knowledge



# Capturing Joint Hierarchy: Project over Semantic Resource



# Capturing Joint Hierarchy: Project over Semantic Resource



# Generalized Loss Function

We optimize a weighted variant of the **ELBO** to train the model:

$$\mathcal{L} = \begin{array}{c} \text{Reconstruct} \\ \text{from sampled } \textit{base} \\ \text{layer frames} \end{array} + \begin{array}{c} \text{Reconstruct} \\ \text{from sampled } \textit{compression} \\ \text{layer frames} \end{array} + \begin{array}{c} \text{Predict the} \\ \text{observed} \\ \text{base} \\ \text{layer frames} \\ \text{(if any)} \end{array} + \begin{array}{c} \text{KL} \\ \text{(Regularizer,} \\ \text{each layer)} \end{array}$$

# Research Questions

1. Is frame inheritance sufficient?
2. How effective are other frame relations?
3. Can our model generate missing events?
4. Can our model generate better event embeddings?



# Datasets (Train)

## Wikipedia Dataset

<i>Partition</i>	<i>#Docs</i>
Train	457k
Dev	16k

- Our Wikipedia dataset is a partial dump of English Wikipedia.
- The FrameNet annotations are automatically extracted.

# Datasets (Test)

- Event Modeling
  - Wikipedia - 21k #docs
- Missing Events
  - Wikipedia - 4k #docs
- Inverse Narrative Cloze
  - Wikipedia Inverse Narrative - 2000 #docs (6 choices)
- Event Similarity Tasks
  - Hard Similarity - 230 event pairs
  - Hard Extension - 1000 event pairs
  - Transitive Sentence Similarity - 108 event pairs

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	—		
SSDVAE	0.9		
ours			
SSDVAE	0.7		
ours			
SSDVAE	0.5		
ours			
SSDVAE	0.4		
ours			
SSDVAE	0.2		
ours			

$\epsilon$ : the average percent of frames observed during training

- Emulate how sufficiently accurate, extractable semantic knowledge may not always be available.
- $\epsilon$  fixed prior to training each model.
- Frames are only observed during training, and never during evaluation

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours			
SSDVAE	0.7	21.19	39.08
ours			
SSDVAE	0.5	31.11	40.18
ours			
SSDVAE	0.4	33.12	47.88
ours			
SSDVAE	0.2	33.31	44.38
ours			

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours		<b>19.39</b>	41.35
SSDVAE	0.7	21.19	39.08
ours			
SSDVAE	0.5	31.11	40.18
ours			
SSDVAE	0.4	33.12	47.88
ours			
SSDVAE	0.2	33.31	44.38
ours		30.15	<b>49.53</b>

✓ SHEM is better able to model longer event sequences

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours		<b>19.39</b>	41.35
SSDVAE	0.7	21.19	39.08
ours			
SSDVAE	0.5	31.11	40.18
ours			
SSDVAE	0.4	33.12	47.88
ours			
SSDVAE	0.2	33.31	44.38
ours		30.15	<b>49.53</b>

✓ SHEMA is better able to model longer event sequences

✓ [See paper] Lexical signal vs. inferred frames: inferred frames and ontological relations are important

✓ [See paper] Hierarchical layer provides useful, less-than-full supervised feedback

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours		<b>19.39</b>	41.35
SSDVAE	0.7	21.19	39.08
ours			
SSDVAE	0.5	31.11	40.18
ours			
SSDVAE	0.4	33.12	47.88
ours			
SSDVAE	0.2	33.31	44.38
ours		30.15	<b>49.53</b>

✓ SHEMA is better able to model longer event sequences

✓ [See paper] Lexical signal vs. inferred frames: inferred frames and ontological relations are important

✓ [See paper] Hierarchical layer provides useful, less-than-full supervised feedback

✓/✗ Able to leverage more observation ( $\epsilon=0.9$ ) or additional structure ( $\epsilon=0.2$ ); ...

# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours		19.39	41.35
SSDVAE	0.7	21.19	39.08
ours		20.26	35.86
SSDVAE	0.5	31.11	40.18
ours		22.16	37.3
SSDVAE	0.4	33.12	47.88
ours		24.02	43.25
SSDVAE	0.2	33.31	44.38
ours		30.15	49.53

✓ SHEMA is better able to model longer event sequences

✓ [See paper] Lexical signal vs. inferred frames: inferred frames and ontological relations are important

✓ [See paper] Hierarchical layer provides useful, less-than-full supervised feedback

✓/✗ Able to leverage more observation ( $\epsilon=0.9$ ) or additional structure ( $\epsilon=0.2$ ); mixed signals  $\Rightarrow$  mixed performance



# 1. Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	–	21.38	24.88
SSDVAE	0.9	19.84	35.56
ours		<b>19.39</b>	41.35
SSDVAE	0.7	21.19	39.08
ours		20.26	35.86
SSDVAE	0.5	31.11	40.18
ours		22.16	37.3
SSDVAE	0.4	33.12	47.88
ours		24.02	43.25
SSDVAE	0.2	33.31	44.38
ours		30.15	<b>49.53</b>

✓ SHEMA is better able to model longer event sequences

✓ [See paper] Lexical signal vs. inferred frames: inferred frames and ontological relations are important

✓ [See paper] Hierarchical layer provides useful, less-than-full supervised feedback

✓/✗ Able to leverage more observation ( $\epsilon=0.9$ ) or additional structure ( $\epsilon=0.2$ ); mixed signals  $\Rightarrow$  mixed performance

☀ Frame inheritance (e.g., IS-A type relations) helpful but not sufficient for hierarchical event modeling

## 2. How effective are other frame relations?

		Low Observation: $\epsilon=0.2$		High Observation: $\epsilon=0.9$	
Model	Frame Relation	Perplexity (↓)	INC Score (↑)	Perplexity (↓)	INC Score (↑)
SSDVAE	-				
SHEM (ours)	Inheritance				
	Using				
	Precedes				
	Causative_of				
	<i>Grouping</i>				
	<i>scenario_only</i>				

## 2. How effective are other frame relations?

		Low Observation: $\epsilon=0.2$		High Observation: $\epsilon=0.9$	
Model	Frame Relation	Perplexity (↓)	INC Score (↑)	Perplexity (↓)	INC Score (↑)
SSDVAE	-	33.31	44.38	19.84	35.56
SHEM (ours)	Inheritance	30.15	49.53	19.39	41.35
	Using	31.37	49.72	19.39	<b>43.23</b>
	Precedes	32.62	47.92	19.57	41.43
	Causative_of	31.82	<b>49.85</b>	19.42	41.38
	<i>Grouping</i>	<b>28.17</b>	48.88	19.44	40.76
	<i>scenario_only</i>	32.01	48.10	<b>18.81</b>	42.29

✓ Existence of broader associations that these relations enable are very helpful

✓ Lower/higher observation consistently better than previous

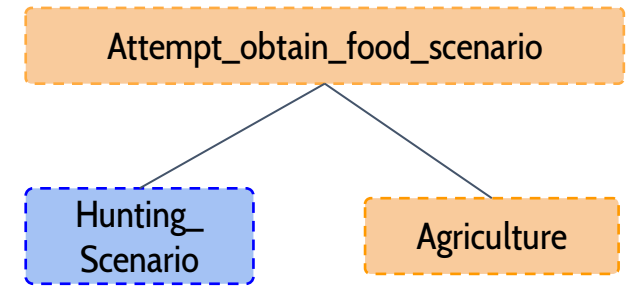
✓ [See paper] even with limited guidance, compression gives valuable feedback

💡 Event modeling could benefit from broader semantic resource coverage: how to encode semantics of any particular relation?

### 3. Can our **SCENARIO** model generate missing events?

Goal: examine the robustness of our model with respect to semantically related missing events in an input sequence

1. Identify sequences where two events have different frames that are contained within the same scenario frame.
- 2.
- 3.



### 3. Can our **SCENARIO** model generate missing events?

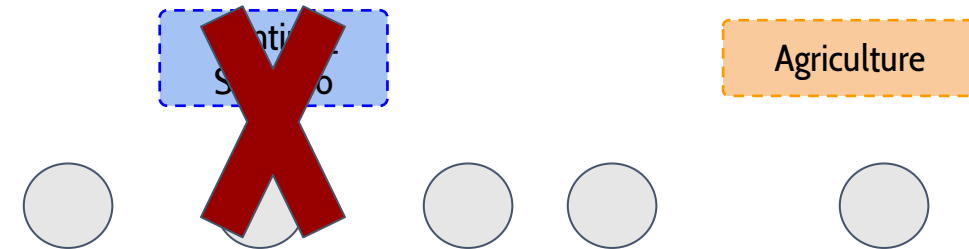
Goal: examine the robustness of our model with respect to semantically related missing events in an input sequence

1. Identify sequences where two events have different frames that are contained within the same scenario frame.
2. Train normally
- 3.

### 3. Can our **SCENARIO** model generate missing events?

Goal: examine the robustness of our model with respect to semantically related missing events in an input sequence

1. Identify sequences where two events have different frames that are contained within the same scenario frame.
2. Train normally
3. To evaluate, consider semantically impoverished input
  - a. **Remove an event** associated with a scenario-connected frame in the input.
  - b. Task: require the model to generate the full, unmodified sequence.



### 3. Are we robust to missing events?

Model	$\epsilon$	Perplexity (Masked Test Data) (↓)
		Base Alone
SSDVAE	0.9	152.44
SHEM: grouping		<u>61.10</u>
SHEM: scenario		<u>63.48</u>
SSDVAE	0.5	182.63
SHEM: grouping		79.74
SHEM: scenario		76.01
SSDVAE	0.2	212.93
SHEM: grouping		89.73
SHEM: scenario		83.86

✓ More observation  $\Rightarrow$   
coarser-grained grouping  
okay

✓ Less observation  $\Rightarrow$   
precise groupings helpful

✓ [See paper] Suggests  
improved capability of SHEM  
to better encapsulate  
abstract meaning of an event  
sequence

## 4. Can our model generate better event embeddings?

- Augment our loss with an event contrastive learning loss, proposed by Gao et al. (2022: SWCC)
- Learn to *align* similar event embeddings, and repel dissimilar ones
- Form event embeddings from decoder & latent compression layer

Model	“Original” Similarity (Acc.)	“Extended” Similarity (Acc.)	Transitivity (Corr.)
SWCC (16)	$78.9 \pm 1.3$	$69.2 \pm 0.9$	<b><math>0.82 \pm 0</math></b>
SWCC (256)	$81.1 \pm 0.4$	$72.6 \pm 1.5$	<b><math>0.82 \pm 0</math></b>
Ours	<b><math>83.3 \pm 2.3</math></b>	<b><math>78.6 \pm 3.0</math></b>	$0.77 \pm 0.04$



# Contrastive Learning vs. Language Modeling

	Similarity (Orig.)	
	SWCC	SHEM
LM only	25.87	67.83
Contrastive only	78.48	67.18
Contrastive + LM	78.91	83.26

Neither contrastive nor LM/MLM loss are as strong as both together

The LM component in our approach is important to overall performance

# Summary

1. We introduce **SHEM**, a novel, hierarchical, semi-supervised event learning model.
2. We show how to use FrameNet for both **observable** event frames and **latent** abstract event frames.
3. More informed signal from compression layer when performing different tasks.
4. Our model can generate better event embeddings for out-of-domain dataset.



[/dipta007/\*\*SHEM\*\*](https://github.com/dipta007/SHEM)

{sroydip1, ferraro}@umbc.edu

# Auxiliary Slides

# 1- Is frame inheritance sufficient?

Model	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE		$21.38 \pm 0.25$	$24.88 \pm 1.35$
SSDVAE	0.9	$19.84 \pm 0.52$	$35.56 \pm 1.70$
ours		<b><math>19.39 \pm 0.30</math></b>	$41.35 \pm 4.25$
SSDVAE	0.7	$21.19 \pm 0.76$	$39.08 \pm 1.55$
ours		$20.26 \pm 1.36$	$35.86 \pm 3.43$
SSDVAE	0.5	$31.11 \pm 0.85$	$40.18 \pm 0.90$
ours		$22.16 \pm 1.62$	$37.30 \pm 3.33$
SSDVAE	0.4	$33.12 \pm 0.54$	$47.88 \pm 3.59$
ours		$24.02 \pm 1.28$	$43.25 \pm 4.97$
SSDVAE	0.2	$33.31 \pm 0.63$	$44.38 \pm 2.10$
ours		$30.15 \pm 2.73$	<b><math>49.53 \pm 1.56</math></b>


## 2- How effective are other frame relations?

Model	Frame Relation	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	-	-	$21.38 \pm 0.25$	$24.88 \pm 1.35$
SSDVAE	-	0.9	$19.84 \pm 0.52$	$35.56 \pm 1.70$
ours	Inheritance		$19.39 \pm 0.53$	$41.35 \pm 4.25$
SSDVAE	Using		$19.39 \pm 0.51$	<b><math>43.23 \pm 2.51</math></b>
ours	Precedes		$19.57 \pm 0.58$	$41.43 \pm 3.02$
SSDVAE	Causative_of		$19.42 \pm 0.57$	$41.38 \pm 2.23$
ours	Inchoative_of		$19.28 \pm 0.32$	$41.35 \pm 3.47$
SSDVAE	Prospective_on		$19.76 \pm 0.97$	$40.53 \pm 2.04$
ours	Subframe		$18.91 \pm 0.15$	$40.35 \pm 2.91$
SSDVAE	<i>Grouping</i>		$19.44 \pm 0.50$	$40.76 \pm 2.86$
ours	<i>scenario_only</i>		<b><math>18.81 \pm 0.50</math></b>	$42.29 \pm 2.86$


## 2- How effective are other frame relations?

Model	Frame Relation	$\epsilon$	Perplexity (↓)	INC Score (↑)
HAQAE	-	-	$21.38 \pm 0.25$	$24.88 \pm 1.35$
SSDVAE	-	<b>0.2</b>	$33.31 \pm 0.63$	$44.38 \pm 2.10$
ours	Inheritance		$30.15 \pm 2.73$	$49.53 \pm 1.56$
SSDVAE	Using		$31.37 \pm 2.08$	$49.72 \pm 1.73$
ours	Precedes		$32.62 \pm 1.65$	$47.92 \pm 2.25$
SSDVAE	Causative_of		$31.82 \pm 3.00$	<b><math>49.85 \pm 0.84</math></b>
ours	Inchoative_of		$32.65 \pm 1.40$	$48.03 \pm 3.35$
SSDVAE	Prospective_on		$33.20 \pm 1.47$	$47.85 \pm 3.53$
ours	Subframe		$32.78 \pm 2.09$	$47.88 \pm 3.31$
SSDVAE	<i>Grouping</i>		<b><math>28.17 \pm 2.26</math></b>	$48.88 \pm 1.37$
ours	<i>scenario_only</i>		$32.01 \pm 0.70$	$48.10 \pm 2.22$

### 3- Can our **GROUPING** model generate missing events?

Model	$\epsilon$	Perplexity (Masked Test Data) (  )		
		Base Alone	Compression Alone	Base + Compression
SSDVAE	0.9	152.44 $\pm$ 3.45	-	-
grouping		<b><u>61.10</u></b> $\pm$ 1.83	94.76 $\pm$ 1.96	76.08 $\pm$ 0.76
SSDVAE	0.7	163.08 $\pm$ 4.52	-	-
grouping		63.50 $\pm$ 3.49	86.23 $\pm$ 0.70	<b><u>73.98</u></b> $\pm$ 2.04
SSDVAE	0.5	182.63 $\pm$ 6.11	-	-
grouping		79.74 $\pm$ 1.79	83.81 $\pm$ 0.96	81.75 $\pm$ 1.13
SSDVAE	0.4	201.55 $\pm$ 4.10	-	-
grouping		84.17 $\pm$ 4.45	81.49 $\pm$ 0.14	82.80 $\pm$ 2.13
SSDVAE	0.2	212.93 $\pm$ 2.54	-	-
grouping		89.73 $\pm$ 4.67	<b><u>77.32</u></b> $\pm$ 0.72	83.28 $\pm$ 2.38

### 3- Can our **SCENARIO** model generate missing events?

Model	$\epsilon$	Perplexity (Masked Test Data) (  )		
		Base Alone	Compression Alone	Base + Compression
SSDVAE	0.9	152.44 $\pm$ 3.45	-	-
scenario		63.48 $\pm$ 4.43	80.94 $\pm$ 7.44	71.60 $\pm$ 4.12
SSDVAE	0.7	163.08 $\pm$ 4.52	-	-
scenario		<b><u>60.06</u></b> $\pm$ 1.68	<b><u>78.36</u></b> $\pm$ 4.52	<b><u>68.58</u></b> $\pm$ 2.30
SSDVAE	0.5	182.63 $\pm$ 6.11	-	-
scenario		76.01 $\pm$ 5.56	78.70 $\pm$ 1.63	77.33 $\pm$ 3.65
SSDVAE	0.4	201.55 $\pm$ 4.10	-	-
scenario		73.77 $\pm$ 7.87	80.00 $\pm$ 1.89	76.77 $\pm$ 4.89
SSDVAE	0.2	212.93 $\pm$ 2.54	-	-
scenario		83.86 $\pm$ 2.74	81.20 $\pm$ 1.17	82.52 $\pm$ 1.93



## 4- Can our model generate better event embeddings?

Model	Batch Size	Hard Similarity (Accuracy %) (↑)		Transitive Score Similarity (↑)
		Original	Extended	
SWCC	16	78.91 ± 1.31	69.20 ± 0.93	0.82 ± 0.00
SWCC	256	81.09 ± 0.43	72.55 ± 1.53	<b>0.82 ± 0.00</b>
ours	16	<b>83.26 ± 2.29</b>	<b>78.63 ± 2.95</b>	0.77 ± 0.04

## 4- Ablation study of ours and SWCC

Model	Training Variant	Hard Similarity (Accuracy %) (↑)		Transitive Score Similarity (↑)
		Original	Extended	
SWCC	Contrastive + LM	78.91 ± 1.31	69.20 ± 0.93	0.82 ± 0.00
	Contrastive Only	78.48 ± 0.83	67.33 ± 0.19	0.78 ± 0.05
	LM only	25.87 ± 1.31	16.78 ± 0.70	0.55 ± 0.04
ours	Contrastive + LM	83.26 ± 2.29	78.63 ± 2.95	0.77 ± 0.04
	Contrastive Only	67.18 ± 1.79	72.75 ± 2.06	0.72 ± 0.02
	LM only	67.83 ± 14.39	62.15 ± 16.52	0.56 ± 0.04