

# WP2: Explainable Neuro-Symbolic Safeguard Framework



The University of Manchester



# Vision - Ambition

To enable:

- **Neural safety through explainability.**
- **Explainability as a support of formal neural verification.**
- End-users to communicate and implement inference controls.
- The creation of safe NN-based applications.



# Vision - Ambition

End-users to  
communicate and  
implement inference  
controls

Neural safety  
through  
explainability

Explainability as an  
enabler of formal neural  
verification.

The creation of safe NN-  
based applications

**WP3**

**WP1**

# Vision - Ambition

## What it means to look inside the black box

### Explainability

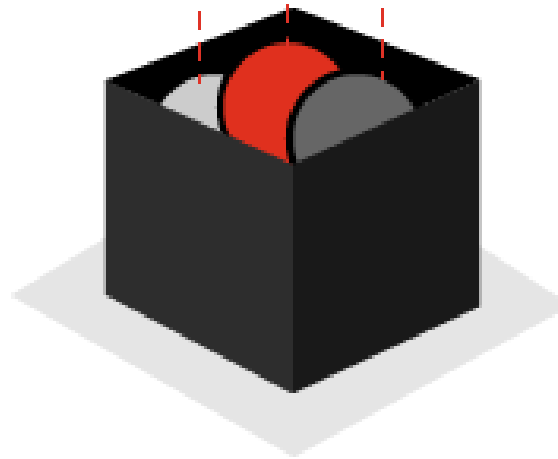
Understanding reasoning  
behind each decision

### Transparency

Understanding of  
AI model decision making

### Provability

Mathematical certainty  
behind decisions



Source: PwC

# Concretely - Safeguards

- Can we ensure that the model does not contradict a medical fact?
- Can we ensure a chatbot is not outputting a racist statement?
- How software developers can program the safeguards into the model?

# More Formally

- **O1**: Develop a novel conceptual/symbolic safeguard mechanism for neuro-symbolic platforms.
- EnnCore will pioneer the use of **neuro-symbolic architectures** and **explainability/interpretability** mechanisms to support end-users specifying a **conceptual safeguard** core to neural-based AI systems.

# Tasks

- **T2.1:** Systematic analysis of neural interpretability methods.
- **T2.2:** Design and implementation of the neurosymbolic safeguard.
- **T2.3:** Evaluation.

Runs from: **M4-M30**



# Neural safety through explainability

Explainability

Semantic controls

Semantic  
probing

Semantic fine-  
tuning

Architectural  
controls

Disentangled  
encodings

Synthetic  
Datasets

Information  
Theoretical Asps.

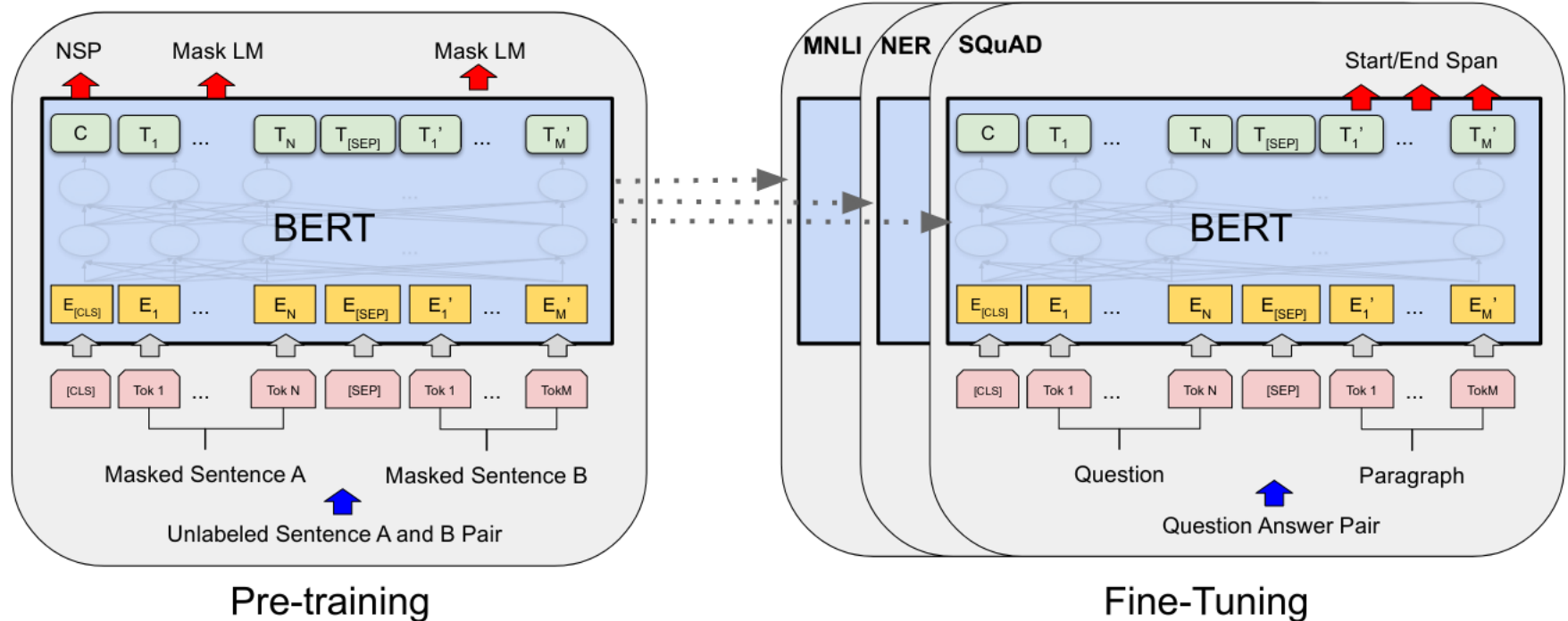
Formal Software  
Verification

Embeddings  
visualisation

Safety Reports

Safety Specs

# Motivational Scenario



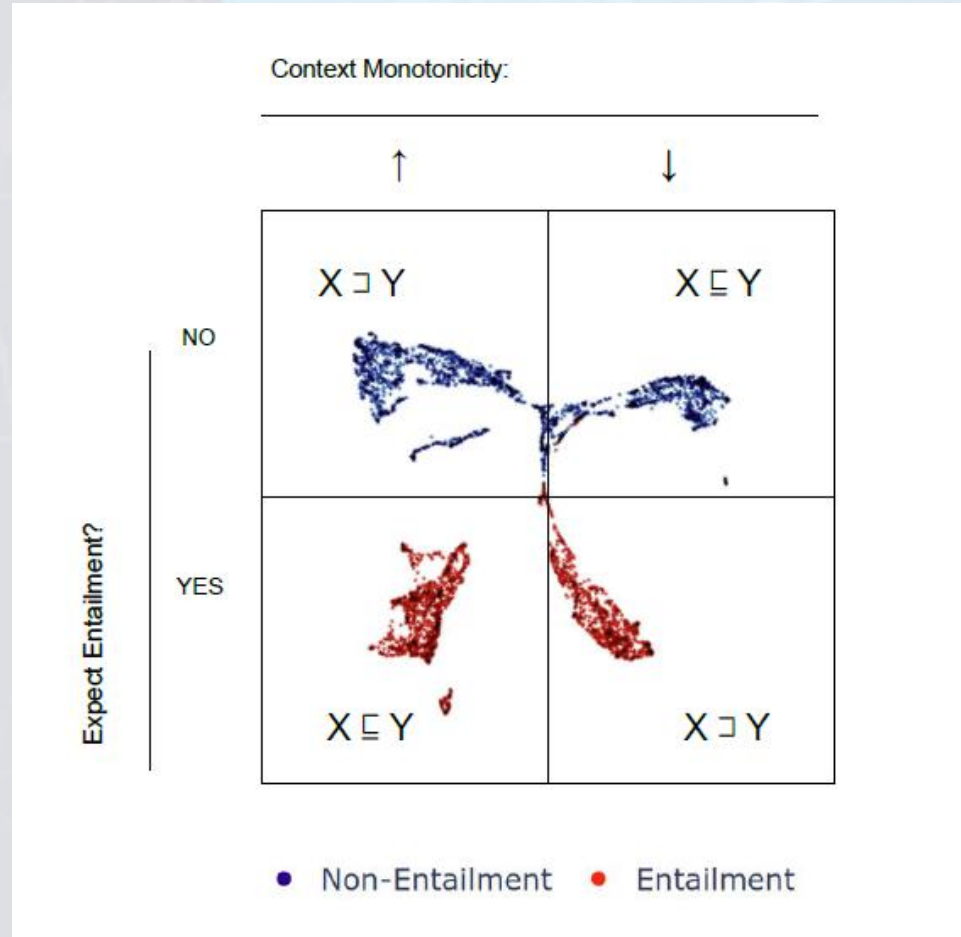
# #1: Development of novel semantic probing/fine-tuning mechanisms

**Monotonicity:** e.g. interpreting negation and generalised quantifiers

		Auxilliary Label
Context	I did not eat any $x$ for breakfast.	↓
Insertion Pair	(fruit, raspberries)	⊃
		NLI Label
Premise	I did not eat any fruit for breakfast.	Entailment
Hypothesis	I did not eat any raspberries for breakfast.	

*Probing Context Monotonicity and Relational Knowledge in Neural Natural Language Inference*

# #1: Development of novel semantic probing/fine-tuning mechanisms



*Probing Context Monotonicity and Relational Knowledge in Neural Natural Language Inference*

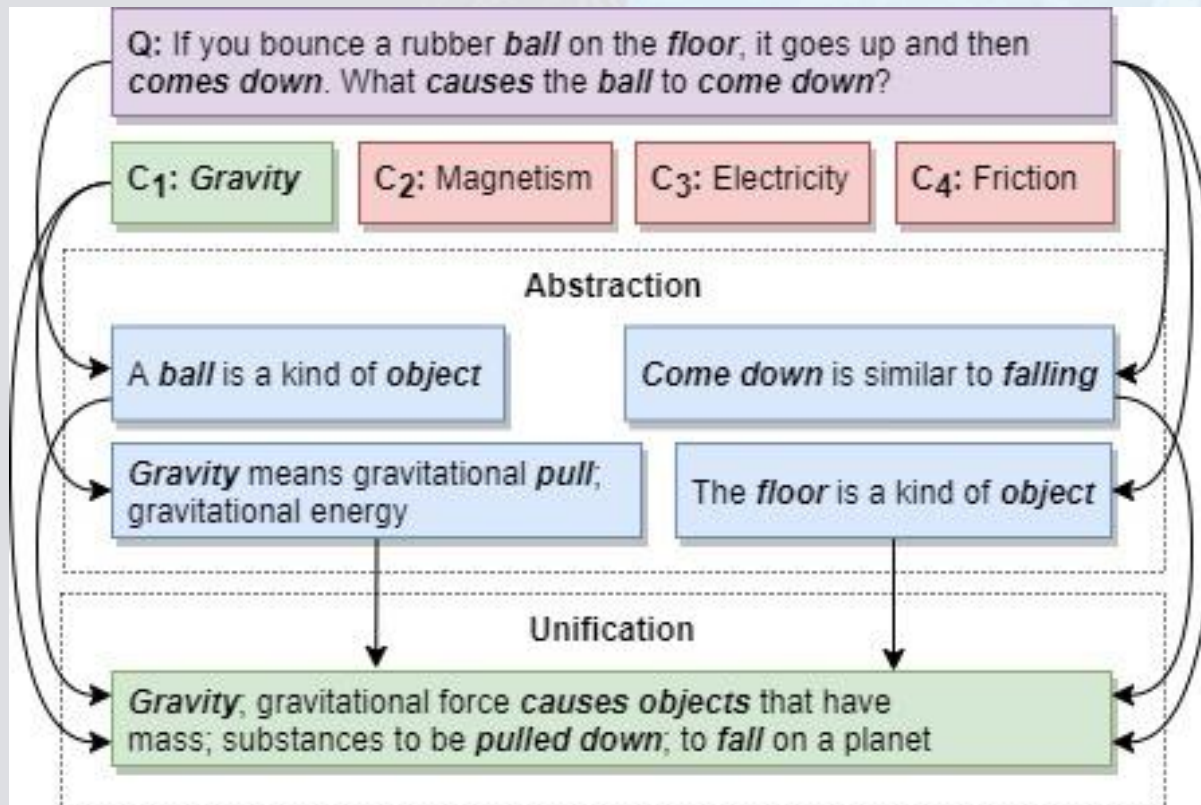
# #1: Development of novel semantic probing/fine-tuning mechanisms

Model	Training Data	Task: Portion of Sequence:	Accuracy at Maximum Selectivity		
			Lexical Relation Classification	Monotonicity Classification	
			XY, concatenated	XY, concatenated	X
Majority Class Baseline	-		0.4050	0.5056	0.5056
Random	-		0.4345	0.5715	0.5715
bert-base	-		0.4605	0.6781	0.7141
bert-large	-		0.4565	0.6766	0.6728
bert-base	SNLI		0.5310	0.7330	0.7281
bert-base	SNLI + HELP		0.5345	0.7455	0.743
facebook/bart-large	MNLI		0.7056	0.7143	0.7513
facebook/bart-large	MNLI + HELP		0.7382	0.7625	0.7643
roberta-large	MNLI		<b>0.8390</b>	0.6771	0.8243
roberta-large	MNLI + HELP		0.7821	<b>0.8742</b>	<b>0.8691</b>

*Probing Context Monotonicity and Relational Knowledge in Neural Natural Language Inference*



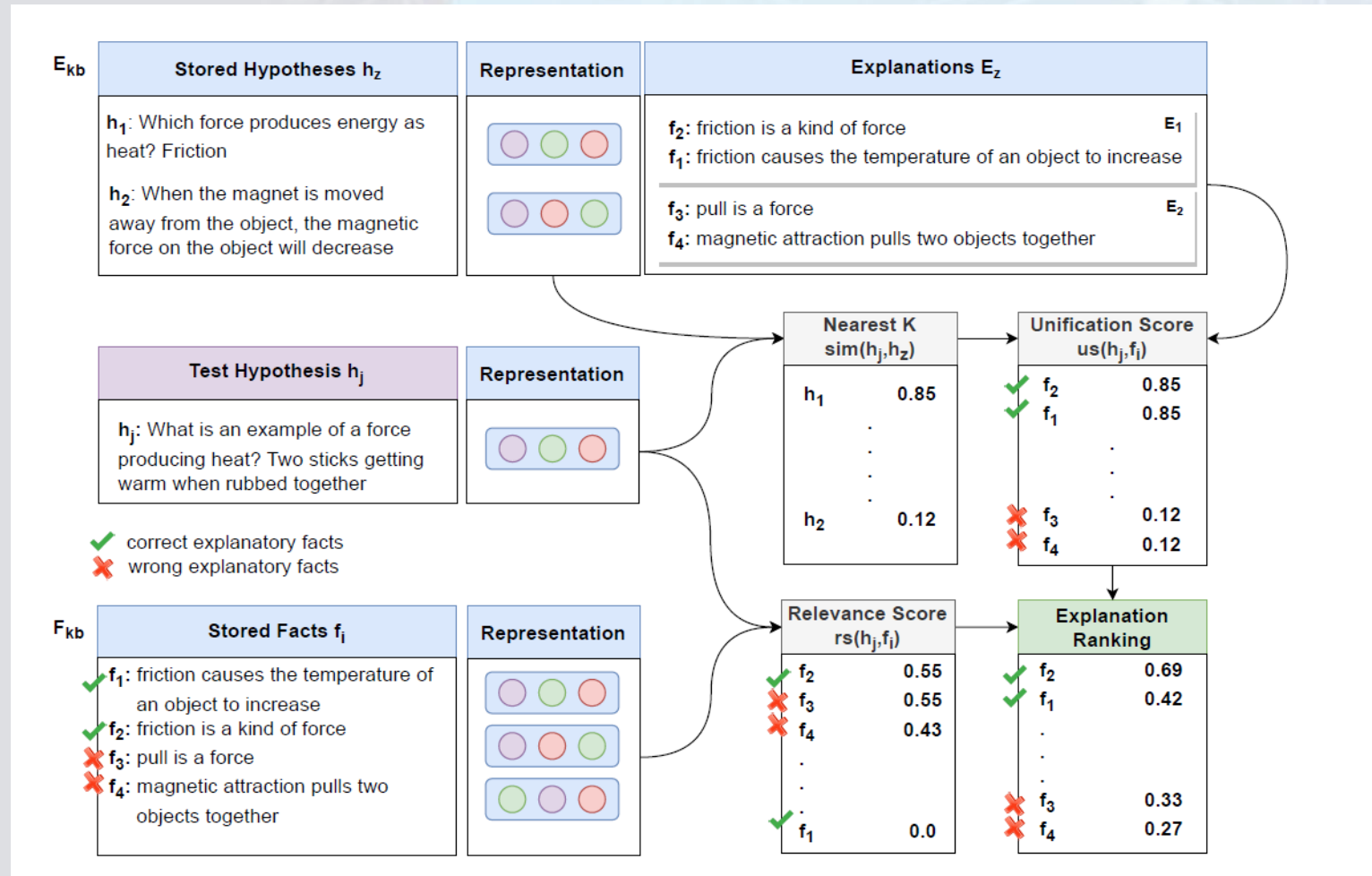
## #2: Semantic controls over complex end-to-end neural architectures



Abstraction

Unification

# #2: Semantic controls over complex end-to-end neural architectures



# #2: Semantic controls over complex end-to-end neural architectures

Question(Q):

What is an example of **force** producing heat

Candidate Answer(C<sub>i</sub>):

Two **sticks** getting warm when **rubbed** together

Grounding Facts:

- [✓] a **stick** is an **object**: F<sub>G1</sub>
- [✓] **friction** is a **force**: F<sub>G2</sub>
- [✗] a **pull** is a **force**: F<sub>G3</sub>
- [✓] to **rub** **together** means to **move against**: F<sub>G4</sub>
- [✗] **rubbing** against something is a kind of **contact**: F<sub>G5</sub>

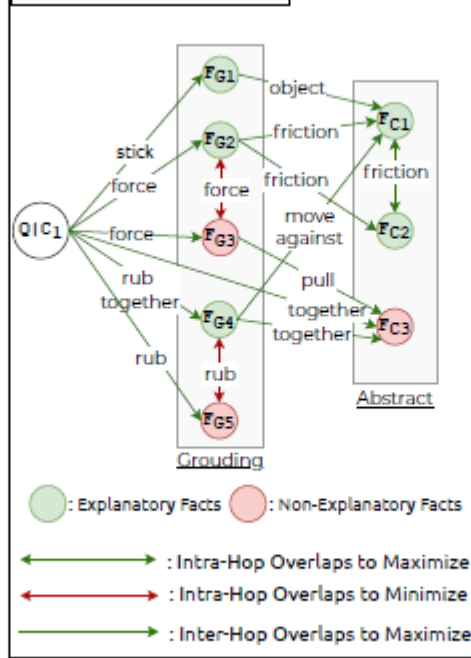
Abstract Facts:

- [✓] **friction** occurs when two **object's** surfaces **move against** each other: F<sub>C1</sub>
- [✓] **friction** causes the temperature of an object to **increases**: F<sub>C2</sub>
- [✗] **magnetic attraction** **pulls** two objects **together**: F<sub>C3</sub>

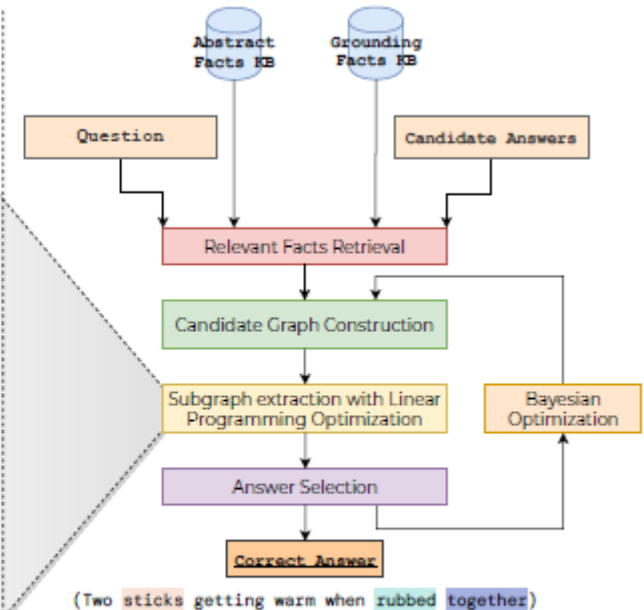
[✓]: Explanatory Facts  
[✗]: Non-Explanatory Facts

(A)

For each Candidate Answer:



(B)



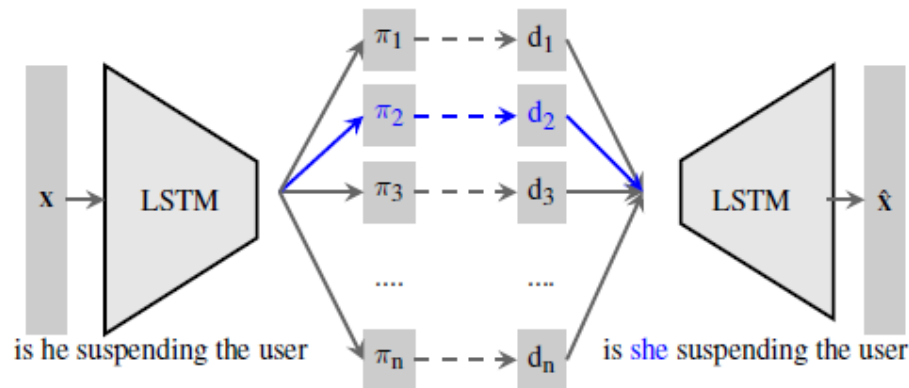
(C)

## #2: Semantic controls over complex end-to-end neural architectures

#	Parameter	Value	
		WT	ARC
1	Question-Abstract overlap ( $\theta_{qa}$ )	0.10	0.09
2	Question-Grounding overlap ( $\theta_{qg}$ )	0.98	0.84
3	Abstract-Abstract overlap ( $\theta_{aa}$ )	0.01	0.11
4	Grounding-Abstract overlap ( $\theta_{ga}$ )	0.14	0.23
5	Grounding-Grounding overlap ( $\theta_{gg}$ )	-0.99	-0.92
6	Abstract Relevance ( $\theta_{ar}$ )	0.03	0.09
7	Grounding Relevance ( $\theta_{gr}$ )	0.36	0.14
8	Edge weight ( $\theta_{ew}$ )	0.80	0.26
9	Node weight ( $\theta_{vw}$ )	0.76	0.67

*ExplanationLP: Abductive Reasoning for Explainable Science Question Answering*

# #3: Development of symbolic disentangled encodings



(a) Model architecture

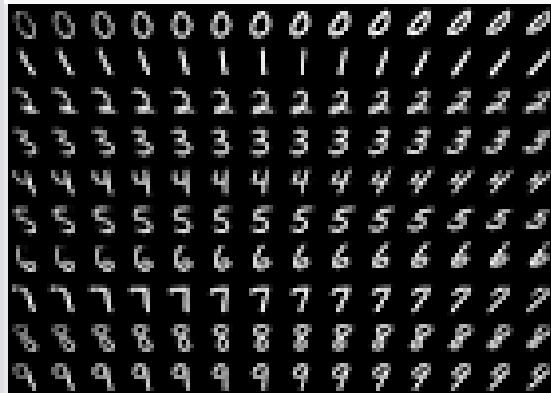
N	Discrete Factor ( $\pi$ )	Sentence example
1	Verb/Object (VO)	is he suspending <b>the computer</b>
2	Gender (G)	is <b>she</b> suspending the user
3	Sentiment (S)	is he <b>not</b> suspending the user
4	Tense (T)	<b>was</b> he suspending the user
5	Subject Number (SN)	are <b>they</b> suspending the user
6	Object Number (ON)	is he suspending the users
7	Question/Answer (QA)	<b>he is</b> suspending the user

(b) Traversal generation

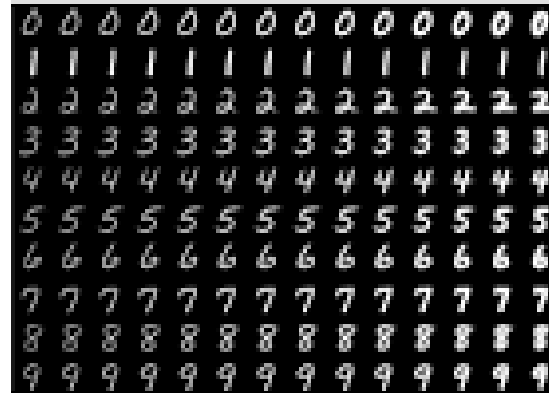


# #3: Development of symbolic disentangled encodings

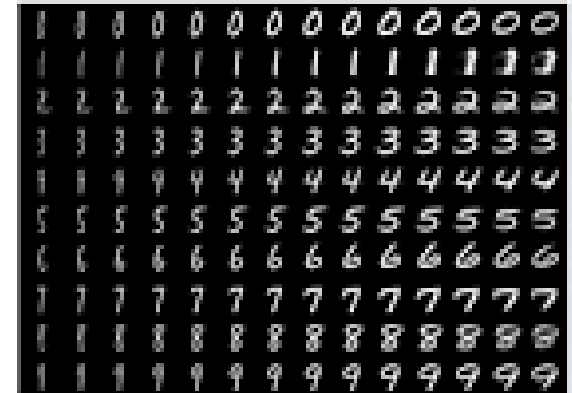
Azimuth



Stroke thickness



Width



# Other aspects

#4: Development of synthetic datasets for textual inference

#5: Connection with feature selection and feature stability

#6: Connection with formal verification

# Discussion points / TODOs

- Initial assumption here is delivering an end-to-end solution
- Closing the gap between different modalities
- Which data modalities to target?
- Planning
- First collaboration points
- Start via synthetic datasets