

# EE-608: Deep Learning For Natural Language Processing: Future Challenges

James Henderson



DLNLP, Lecture 12

# Announcements

- ▶ Project presentations next week, during both lecture and exercises sessions. Let us know soon if you have time constraints.
- ▶ We encourage everyone to do their presentations in-person, but doing it by Zoom is allowed, including hybrid presentations. Please organise this with your team members.
- ▶ Watching the presentations of other groups via Zoom will be possible.

# Outline

Understanding Transformers as Structure Induction Models

Challenges: Abstraction

Challenges: Compositionality

Challenges: Causal Latent Variables

Summary of Future Challenges

# Outline

## Understanding Transformers as Structure Induction Models

Challenges: Abstraction

Challenges: Compositionality

Challenges: Causal Latent Variables

Summary of Future Challenges

# The Unstoppable Rise of Computational Linguistics in Deep Learning

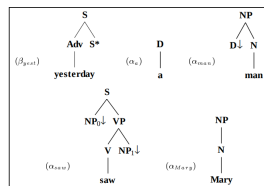
This section is from the ACL 2020 theme paper:

James Henderson. 2020. “The Unstoppable Rise of Computational Linguistics in Deep Learning”. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6294–6306, Online.

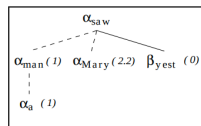
# The Nature of Language

Grammar Formalisms are designed to capture the generalisations found in natural languages.

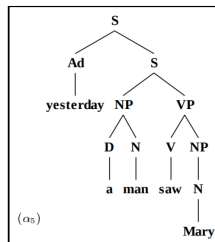
- ▶ Multiple levels of representation
- ▶ each with multiple entities
- ▶ each with multiple properties
- ▶ and relations between entities
- ▶ Rules specify the allowable structures
- ▶ Grammar formalisms specify the allowable rules



TAG Grammar



Derivation structure

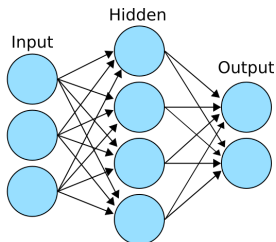
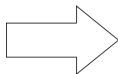


Linguistic structure

# Connectionism

Brains have distributed representations and learning.

- ▶ Vector spaces formalise distributed representations
- ▶ Backpropagation learning is astoundingly effective
- ▶ Multi-layer and recurrent neural networks are powerful approximators

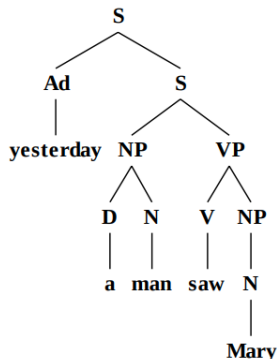


# Variable Binding



Questioning the adequacy of vector spaces.

- ▶ Our representations segment the world into entities
- ▶ Variable binding: How are multiple entities represented in the brain?
- ▶ Systematicity: How can we learn rules which generalise across entities?



# Overview

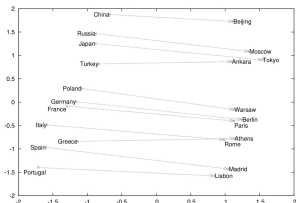
Our understanding of the nature of language from computational linguistics has fundamentally influenced deep learning architectures.

1. Inducing Features of Entities
2. Inducing Relations between Entities
3. Inducing Entities?

# Early Neural Network Models of Language

Neural networks were successfully integrated in classical natural language understanding models

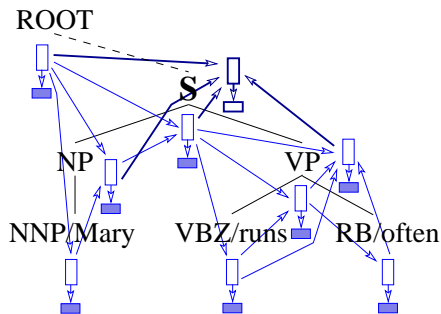
- ▶ Learned vectors replace categorical labels
  - ▶ Word embeddings
  - ▶ Contextualised token embeddings
  - ▶ Parse history embeddings
- ▶ Neural probability estimation replaces feature engineering
  - ▶ Subtree scores
  - ▶ Parser action scores



# Modelling Derivation Structures

But successful neural approaches kept the same structures

- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



# Modelling Derivation Structures

But successful neural approaches kept the same structures

- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure

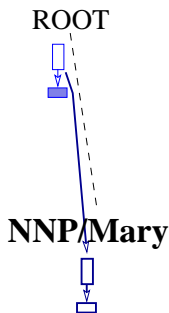
**ROOT**



# Modelling Derivation Structures

But successful neural approaches kept the same structures

- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



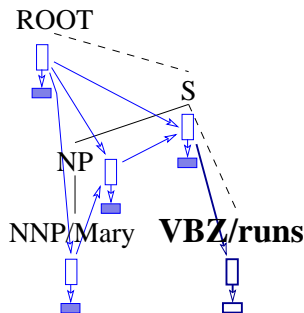




# Modelling Derivation Structures

But successful neural approaches kept the same structures

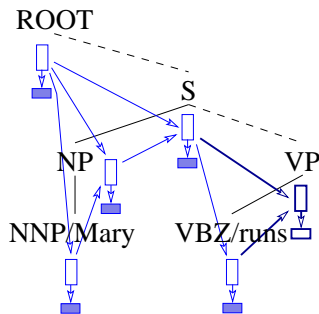
- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



# Modelling Derivation Structures

But successful neural approaches kept the same structures

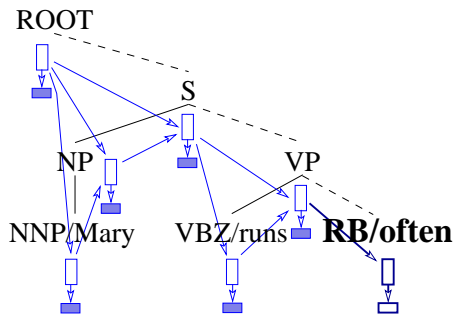
- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



# Modelling Derivation Structures

But successful neural approaches kept the same structures

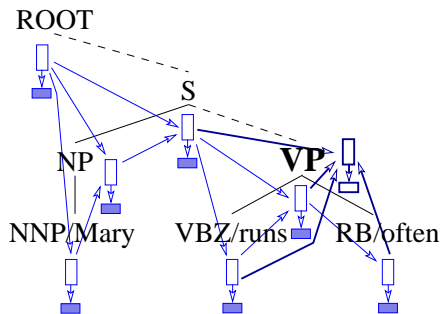
- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



# Modelling Derivation Structures

But successful neural approaches kept the same structures

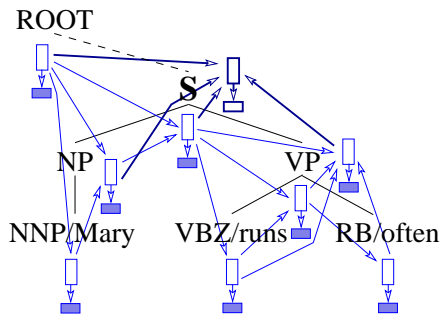
- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



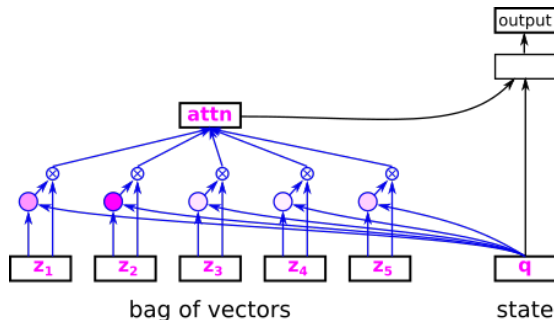
# Modelling Derivation Structures

But successful neural approaches kept the same structures

- ▶ Linguistic structures reflect regularities
- ▶ Derivation structures reflect the linguistic structure
- ▶ Hand-coded neural network model structures reflect the derivation structure



# Attention Induces Structure



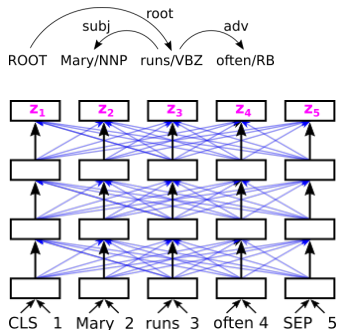
Attention-based models are induced structure models, not sequence models

- ▶ Attention weights determine the model structure
- ▶ Model structure reflects linguistic structure
- ▶ Attention treats a sequence as a bag-of-vectors

# Transformers and Variable Binding

Transformers have variable binding.

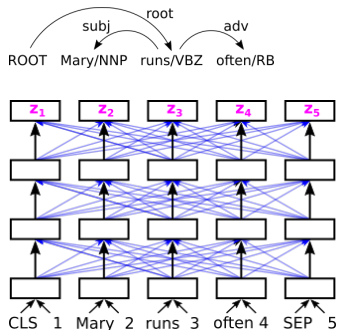
- ▶ Each “position” is an entity
- ▶ Each vector embeds the properties of its entity
- ▶ Each pair of vectors embeds the relations between their entities



# Transformers and Systematicity

Transformers have systematicity.

- ▶ Rules are parameterised in terms of an entity's properties and relations
- ▶ Parameters are shared across entities
- ▶ Rules learned for one entity inherently generalise to other entities



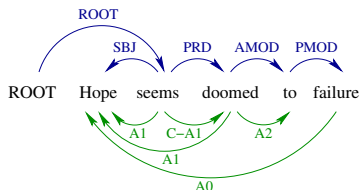
# Future Directions

What is left to learn?

- ▶ How do we learn the set of entities?



- ▶ Currently, the set of positions are pre-defined, with hand-coding (e.g. words), or preprocessing (e.g. BPE)
- ▶ How do we learn the set of levels?
  - ▶ Linguistic levels seem to be language universal
  - ▶ Are they innate?
  - ▶ Can they be learned at all?



## Summary

Our understanding of the nature of language from computational linguistics has fundamentally influenced deep learning architectures.

- ▶ Vector spaces and sequence-to-sequence models are **not adequate**
- ▶ Attention-based models **induce structured representations** over entities
- ▶ Transformers use bag-of-vector representations (variable binding)
- ▶ Transformers inherently generalise across entities (systematicity)

But levels of representation and their entities are predefined.  
Can they be learned from data?

# Outline

Understanding Transformers as Structure Induction Models

**Challenges: Abstraction**

Challenges: Compositionality

Challenges: Causal Latent Variables

Summary of Future Challenges

## Limitations of current LLMs

- LLMs excel at many tasks but still struggle on low-resource languages, systematic generalisation, and abstraction.
- Need: curated diagnostic tasks that probe representation of objects, structure-dependency and systematicity.

## Do LLMs know the grammar and the meaning of the language they seem to use and generate so well?

The truck was loaded with hay. The hay was loaded onto the truck.  
Locative Instrumental Instrumental Locative

- Verb alternations require object identification, attribute assignment and systematic correspondence.
- Can current LLMs perform these linguistic operations?

## More problems and underlying rules

**Examples of verb alternations** (Levin 1993) (Italian, English, Hebrew)

The man rolled the dice/ The dice rolled (manner of motion)

The cook burnt the cake/ The cake burnt (Change of state)

The artist paints the portrait/ The artist paints (object-drop)

**Example of agreement** (French, Italian, English, Rumanian)

L'ordinateur avec le programme est en panne.

Les ordinateurs avec le programme sont en panne.

L'ordinateur avec le programme de l'expérience est en panne.

Les ordinateurs avec le programme de l'expérience sont en panne.

## The architectural and data-centric approach

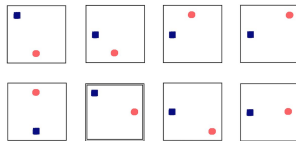
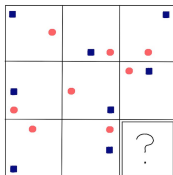
- We develop curated synthetic data on a large scale. We develop **data and diagnostic probes** that help us understand their current generalisation abilities.

Merlo 2023 Arxiv

An et al 2023 EACL, Samo et al 2023, Findings EMNLP, Merlo 2023,  
Findings Emnlp,

# The task: BLMs

RPMs



BLMs

|   | CONTEXT |        |        |       |
|---|---------|--------|--------|-------|
| 1 | NP-sg   | PP1-sg | VP-sg  |       |
| 2 | NP-pl   | PP1-sg | VP-pl  |       |
| 3 | NP-sg   | PP1-pl | VP-sg  |       |
| 4 | NP-pl   | PP1-pl | VP-pl  |       |
| 5 | NP-sg   | PP1-sg | PP2-sg | VP-sg |
| 6 | NP-pl   | PP1-sg | PP2-sg | VP-pl |
| 7 | NP-sg   | PP1-pl | PP2-sg | VP-sg |
| 8 | ???     |        |        |       |

|   | ANSWERS |        |           |       |             |
|---|---------|--------|-----------|-------|-------------|
| 1 | NP-pl   | PP1-pl | PP2-sg    | VP-pl | <b>Corr</b> |
| 2 | NP-pl   | PP1-pl | et PP2-sg | VP-pl | Coord       |
| 3 | NP-pl   | PP1-pl |           | VP-pl | WNA         |
| 4 | NP-pl   | PP1-sg | PP1-sg    | VP-pl | WN1         |
| 5 | NP-pl   | PP1-pl | PP2-pl    | VP-pl | WN2         |
| 6 | NP-pl   | PP1-pl | PP2-pl    | VP-sg | AEV         |
| 7 | NP-pl   | PP1-sg | PP2-pl    | VP-sg | AEN1        |
| 8 | NP-pl   | PP1-pl | PP2-sg    | VP-sg | AEN2        |

## The task: BLMs

|   | CONTEXT |        |        |       |
|---|---------|--------|--------|-------|
| 1 | NP-sg   | PP1-sg | VP-sg  |       |
| 2 | NP-pl   | PP1-sg | VP-pl  |       |
| 3 | NP-sg   | PP1-pl | VP-sg  |       |
| 4 | NP-pl   | PP1-pl | VP-pl  |       |
| 5 | NP-sg   | PP1-sg | PP2-sg | VP-sg |
| 6 | NP-pl   | PP1-sg | PP2-sg | VP-pl |
| 7 | NP-sg   | PP1-pl | PP2-sg | VP-sg |
| 8 | ???     |        |        |       |

|   | ANSWERS |        |           |       |       |
|---|---------|--------|-----------|-------|-------|
| 1 | NP-pl   | PP1-pl | PP2-sg    | VP-pl | Corr  |
| 2 | NP-pl   | PP1-pl | et PP2-sg | VP-pl | Coord |
| 3 | NP-pl   | PP1-pl |           | VP-pl | WNA   |
| 4 | NP-pl   | PP1-sg | PP1-sg    | VP-pl | WN1   |
| 5 | NP-pl   | PP1-pl | PP2-pl    | VP-pl | WN2   |
| 6 | NP-pl   | PP1-pl | PP2-pl    | VP-sg | AEV   |
| 7 | NP-pl   | PP1-sg | PP2-pl    | VP-sg | AEN1  |
| 8 | NP-pl   | PP1-pl | PP2-sg    | VP-sg | AEN2  |

What kind of linguistic objects the LLM actually learn and manipulate?

Are they tokens?

Does the LLM discover something more linguistically justified like phrases and constituents?

Are the inner representations of NNs information about chunks mapped to semantic roles? Are long-distance dependencies in agreement represented?

Are the objects and attributes represented by higher-level abstractions or based on low-level cues?

Do the abstractions hold across languages?

# Blackbird Language Matrices: A Framework to Investigate the Linguistic Competence of LLMs

---

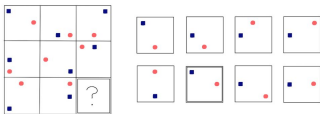
Paola Merlo

Computational Learning and Computational Linguistics group

[www.idiap.ch](http://www.idiap.ch)

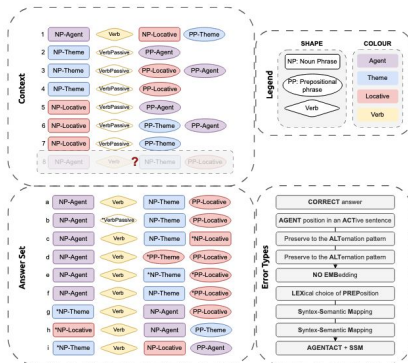
*Slide from Paola Merlo*

# The BLM task



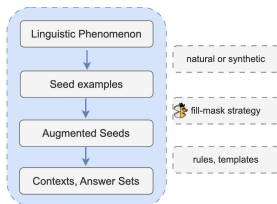
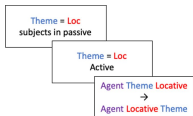
Beyond performance: human-like generalisation in NN

- Rule-based generalisation
- With little data
- Across languages



# Creating matrices and datasets

| CONTEXT |          |          |                 |
|---------|----------|----------|-----------------|
| 1       | NP-Agent | Verb     | NP-Loc PP-Theme |
| 2       | NP-Theme | VerbPass | PP-Agent        |
| 3       | NP-Theme | VerbPass | PP-Loc PP-Agent |
| 4       | NP-Theme | VerbPass | PP-Loc          |
| 5       | NP-Loc   | VerbPass | PP-Agent        |
| 6       | NP-Loc   | VerbPass | PP-Theme        |
| 7       | NP-Loc   | VerbPass | PP-Theme        |
| 78      | NP-Agent | Verb     | NP-Theme PP-Loc |



## Agreement

English French  
Italian Romanian

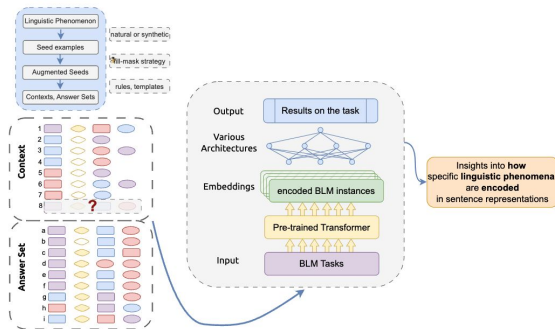
## Spray/Load

She loads the truck with hay  
She loads hay into the truck

## Causatives

Mary opened the door  
The door opened.

# The framework pipeline



V. Nastase and P. Merlo. (2024). *Tracking linguistic information in transformer-based sentence embeddings through targeted sparsification*. In: Proceedings of the 9th Workshop on Representation Learning for NLP (Repl4NLP-2024), pages 203–214, Bangkok, Thailand.

Slide from Paola Merlo

# Experiments

Can BLMs be solved by current models? At what accuracy?

(benchmarks)

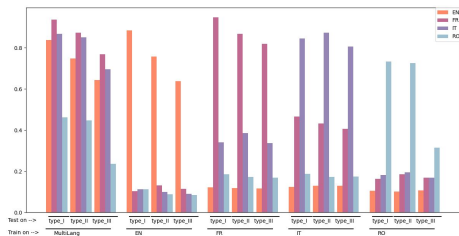
Do sentence embeddings encode chunk/constituent structure?

Can compressed representations reveal systematic structure across sentences?

# The linguistic results

No cross-lingual abstraction

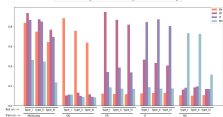
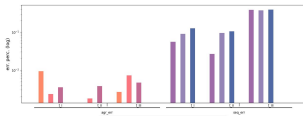
Nastase et al 2024, Clic-it.



Slide from Paola Merlo

# Results

| Baseline-FFNN    |                |                | VAE_5_1x2 |                |                |      |
|------------------|----------------|----------------|-----------|----------------|----------------|------|
|                  | test on type_1 | test on type_2 |           | test on type_1 | test on type_2 |      |
| train on type_1  | 0.98           | 0.77           | 0.69      | 1              | 0.83           | 0.78 |
| train on type_2  | 0.96           | 0.79           | 0.72      | 0.99           | 0.9            | 0.85 |
| train on type_20 | 0.91           | 0.77           | 0.76      | 0.98           | 0.88           | 0.87 |



The task

Good overall performance, sensitive to lexicalisation

Error analysis

Good on language, less well on sequence reasoning

Across languages

No cross-lingual abstraction.

# Outline

Understanding Transformers as Structure Induction Models

Challenges: Abstraction

**Challenges: Compositionality**

Challenges: Causal Latent Variables

Summary of Future Challenges

## Compositional Representations and Systematic Generalization

- Systematicity: The ability to produce/understand some sentences is intrinsically connected to the ability to produce / understand certain others. This means there is a “definite and predictable pattern among the sentences we understand”
- E.g. any speaker that understands the sentence “John loves Mary” should be able to understand “Mary loves John”.

# Compositional Representations and Systematic Generalization

Compositionality: closely related to the idea of systematicity is the principle of compositionality.

## Rough Definition:

“The meaning of an expression is a function of the meaning of its parts”



## More concrete definition (Montague):

A homomorphism from syntax (structure) to semantics (meaning). That is, meaning of the whole is a function of immediate constituents (as determined by syntax)

## Are human languages compositional?

Brown Cow = *Brown objects*  $\cap$  *Cows*



Red Rabbit = *Red objects*  $\cap$  *Rabbits*



Kicked the Ball = *Kicked(Ball, Agent)*



## Are human languages compositional?

Brown Cow = *Brown objects*  $\cap$  *Cows*



Red Rabbit = *Red objects*  $\cap$  *Rabbits*



Kicked the Ball = *Kicked(Ball, Agent)*



Red Herring  $\neq$  *Red things*  $\cap$  *Herring*



Kicked the bucket  $\neq$  *Kicked(Bucket, Agent)*

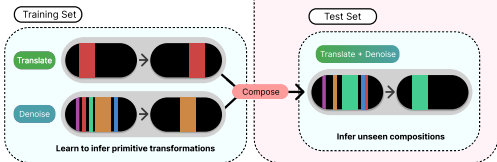


# Compositional Generalisation Tasks

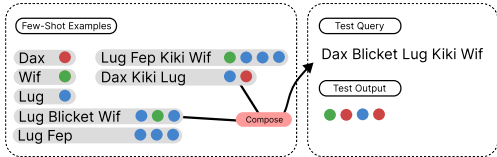
## A. 1-D ARC



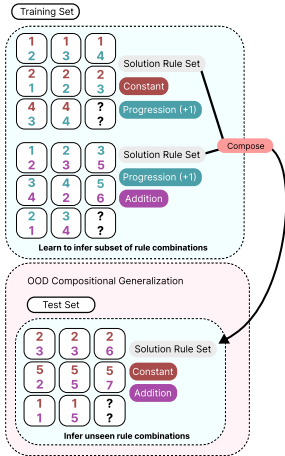
## B. Composition in ARC



## D. Linguistic Systematicity Task

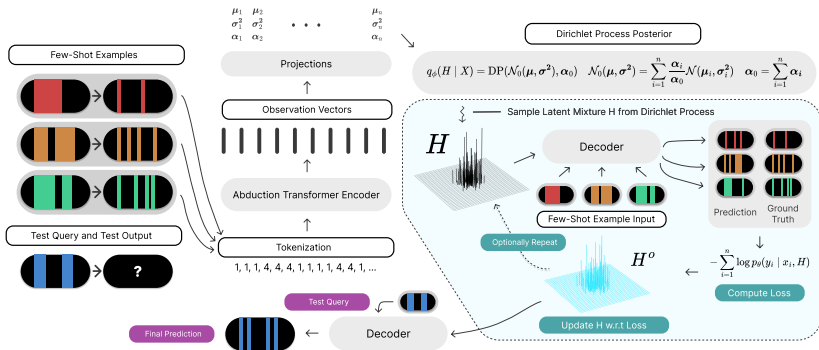


## C. Raven's Progressive Matrices



Test sets contain unseen compositions of seen primitives.

# Abduction Transformer Model



Search for a latent program which solves the few-shot examples.

- ▶ An encoder tries to guess the latent program
- ▶ Gradient search fine-tunes the latent program
- ▶ Uses a nonparametric latent space to generalise from simple to complex programs

# Outline

Understanding Transformers as Structure Induction Models

Challenges: Abstraction

Challenges: Compositionality

**Challenges: Causal Latent Variables**

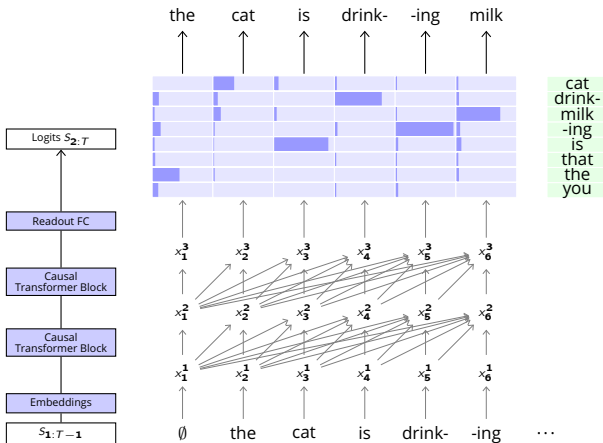
Summary of Future Challenges

Selected slides from (arXiv 2025):

# The Free Transformer

François Fleuret





Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O:

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, 0: physics. I: the square root of two is irrational, 0: mathematics. I: the set of prime numbers is infinite, 0: mathematics. I: gravity is proportional to the mass, 0: [physics](#).

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O:

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O: [mathematics](#).

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O: [mathematics](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas are awesome, O:

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O: [mathematics](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas are awesome, O: [positive](#).

Such a model, trained only to generate text, can already be put to use for instance for classification.

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O: [mathematics](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas are awesome, O: [positive](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas taste like cardboard, O:

Such a model, trained only to generate text, can already be put to use for instance for classification.

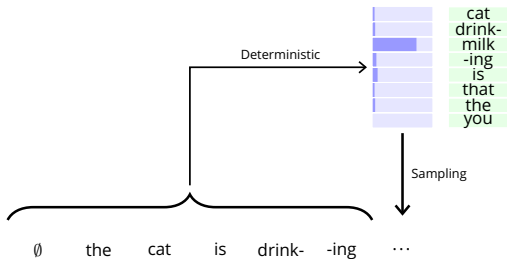
I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: gravity is proportional to the mass, O: [physics](#).

I: water boils at 100 degrees, O: physics. I: the square root of two is irrational, O: mathematics. I: the set of prime numbers is infinite, O: mathematics. I: squares are rectangles, O: [mathematics](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas are awesome, O: [positive](#).

I: I love apples, O: positive. I: music is my passion, O: positive. I: my job is boring, O: negative. I: frozen pizzas taste like cardboard, O: [negative](#).

## The problem with autoregression



The only sampling in a Transformer are the tokens.

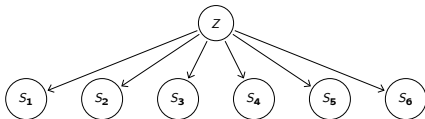
Consider a corpus of positive and negative online reviews.

$$Z \sim \mathcal{U}(\{-1, 1\}), S \sim \mu_Z.$$

An AR model has no way of implementing this factorized distribution, it cannot “decide” beforehand what type of review to generate.

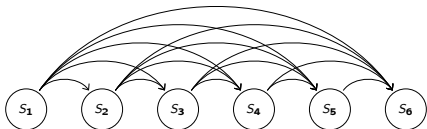
It would estimate on the fly the negativity / positivity of what it has written so far to be consistent.

The “true model” of a joint distribution, with latents, is usually simpler and more robust than the “latentless” AR model.



$$P(S_{t+1} = 1 \mid Z, S_1, \dots, S_t) = (1 - \epsilon)Z + \epsilon(1 - Z)$$

The “true model” of a joint distribution, with latents, is usually simpler and more robust than the “latentless” AR model.



$$P(X_{t+1} = 1 \mid X_1 = x_1, \dots, X_t = x_t) = \frac{\left(\frac{\epsilon}{1-\epsilon}\right)^{\sum_{s=1}^t x_s} (1-\epsilon)^t \epsilon + \left(\frac{1-\epsilon}{\epsilon}\right)^{\sum_{s=1}^t x_s} \epsilon^t (1-\epsilon)}{\left(\frac{\epsilon}{1-\epsilon}\right)^{\sum_{s=1}^t x_s} (1-\epsilon)^t + \left(\frac{1-\epsilon}{\epsilon}\right)^{\sum_{s=1}^t x_s} \epsilon^t}.$$

We could prefix every training sample with a token  $Z$  indicating if the review is positive or not.

We could prefix every training sample with a token  $Z$  indicating if the review is positive or not.

Reasoning post-training tries to address it unsupervised, however:

- + it requires a trained model,
- + current approaches cast it as reinforcement learning,
- + the conditioning variable are discrete tokens,
- + it cannot deal with stochastic responses.

We propose instead to let the model build latent variables

$$Y_r = f_r(S_1, \dots, S_t, Y_1, \dots, Y_{r-1}, Z_r; \theta)$$

where  $Z_r$  is sampled from a random generator. This can be interpreted as making decisions beside the token choices.

We propose instead to let the model build latent variables

$$Y_r = f_r(S_1, \dots, S_t, Y_1, \dots, Y_{r-1}, Z_r; \theta)$$

where  $Z_r$  is sampled from a random generator. This can be interpreted as making decisions beside the token choices.

Sampling from such a trained model is trivial: sample  $z$ , and run the AR process as usual to sample  $P_\theta(S | Z = z)$ .

Training, however, is far more involved.

# The Variational Autoencoder

Given a training sample  $s$ , we want to maximize

$$P_{\theta}(S = s) = \mathbb{E}_{z \sim P(Z)} \left[ P_{\theta}(S = s \mid Z = z) \right]$$

Given a training sample  $s$ , we want to maximize

$$P_{\theta}(S = s) = \mathbb{E}_{z \sim P(Z)} \left[ P_{\theta}(S = s \mid Z = z) \right]$$

Since  $P_{\theta}(S = s \mid Z = z)$  is a full-fledged model, there is no closed form for  $P_{\theta}(S = s)$ , and numerical integration requires “good”  $Z$ s.

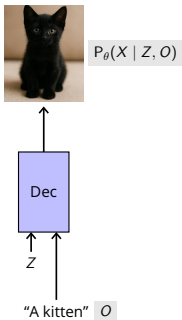
Given a training sample  $s$ , we want to maximize

$$P_{\theta}(S = s) = \mathbb{E}_{Z \sim P(Z)} \left[ P_{\theta}(S = s \mid Z = z) \right]$$

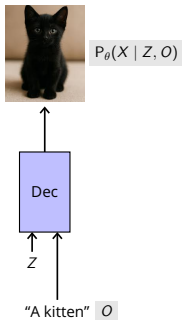
Since  $P_{\theta}(S = s \mid Z = z)$  is a full-fledged model, there is no closed form for  $P_{\theta}(S = s)$ , and numerical integration requires “good”  $Z$ s.

The Variational Autoencoder proposed by Kingma and Welling (2013) provides a formal derivation to train a sampler  $q_{\theta}(Z; S = s)$ .

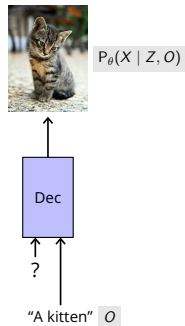
## Inference



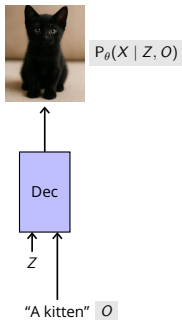
## Inference



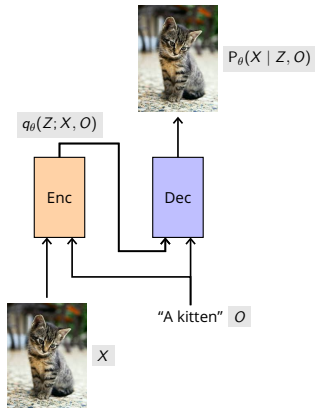
## Training



## Inference



## Training



$$\log P_{\theta}(S = s)$$

$$\log P_{\theta}(S = s)$$

$$\geq \log P_{\theta}(S = s) - \overbrace{\mathbb{D}_{\text{KL}}(q_{\theta}(Z; S = s) \parallel P_{\theta}(Z | S = s))}^{\text{How bad the sampler is}}$$

$$\begin{aligned}
& \log P_{\theta}(S = s) \\
& \geq \log P_{\theta}(S = s) - \overbrace{\mathbb{D}_{\text{KL}}\left(q_{\theta}(Z; S = s) \parallel P_{\theta}(Z \mid S = s)\right)}^{\text{How bad the sampler is}} \\
& = \underbrace{\mathbb{E}_{z \sim q_{\theta}(Z; S = s)} \left[ \log \frac{P_{\theta}(S = s, Z = z)}{q_{\theta}(Z = z; S = s)} \right]}_{\text{"Evidence Lower Bound" (aka ELBO)}}
\end{aligned}$$

$$\begin{aligned}
& \log P_{\theta}(S = s) \\
& \geq \log P_{\theta}(S = s) - \overbrace{\mathbb{D}_{\text{KL}}\left(q_{\theta}(Z; S = s) \parallel P_{\theta}(Z | S = s)\right)}^{\text{How bad the sampler is}} \\
& = \underbrace{\mathbb{E}_{z \sim q_{\theta}(Z; S = s)} \left[ \log \frac{P_{\theta}(S = s, Z = z)}{q_{\theta}(Z = z; S = s)} \right]}_{\text{"Evidence Lower Bound" (aka ELBO)}} \\
& = \underbrace{\mathbb{E}_{z \sim q_{\theta}(Z; S = s)} \left[ \log P_{\theta}(S = s | Z = z) \right]}_{\text{-cross-entropy}} - \underbrace{\mathbb{D}_{\text{KL}}\left(q_{\theta}(Z; S = s) \parallel P(Z)\right)}_{\text{How much we cheat}} \\
& \underbrace{\hspace{10em}}_{\text{The quantity we maximize during training}}
\end{aligned}$$

The model  $q_{\theta}(Z = z; S = s)$  can be envisioned as an encoder that extracts the necessary information from  $s$ .

E.g. “ $s$  is a negative review,  $\theta$  should be improved for  $z = -1$ ”.

The model  $q_{\theta}(Z = z; S = s)$  can be envisioned as an encoder that extracts the necessary information from  $s$ .

E.g. “ $s$  is a negative review,  $\theta$  should be improved for  $z = -1$ ”.

The term

$$\mathbb{D}_{\text{KL}}\left(q_{\theta}(Z; S = s) \parallel P(Z)\right)$$

reflects how much information about  $s$  the encoder provides to the decoder, and must be controlled carefully during training.

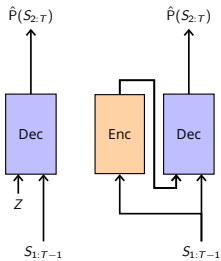
# The Free Transformer



Decoder Transformer



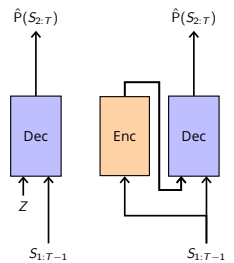
Decoder Transformer



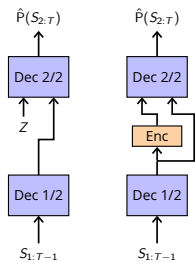
VAE + Transformer



Decoder Transformer



VAE + Transformer



Free Transformer

We use for  $Z$  a one-hot vector of dimension  $2^{16}$ , comparable to the vocabulary size of  $2^{17}$

We use for  $Z$  a one-hot vector of dimension  $2^{16}$ , comparable to the vocabulary size of  $2^{17}$

Given  $L \in \mathbb{R}^H$ , the “Binary Mapper” interprets them as bits logits and

- + Computes the distribution  $p$  over  $\{0, \dots, 2^H - 1\}$ .
- + Samples  $K \sim p$ .
- + Outputs  $Y_d = \delta_{d=K} + p_d - \text{detach}(p_d)$ .

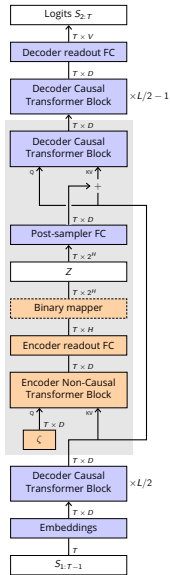
We use for  $Z$  a one-hot vector of dimension  $2^{16}$ , comparable to the vocabulary size of  $2^{17}$

Given  $L \in \mathbb{R}^H$ , the “Binary Mapper” interprets them as bits logits and

- + Computes the distribution  $p$  over  $\{0, \dots, 2^H - 1\}$ .
- + Samples  $K \sim p$ .
- + Outputs  $Y_d = \delta_{d=K} + p_d - \text{detach}(p_d)$ .

The vector  $L$  is:

- + Constant zero for  $P(Z)$ .
- + The output of the encoder for  $q_\theta(Z; S = s)$ .



The KL divergence can be expressed per token

$$\mathbb{D}_{\text{KL}}\left(q(Z_t; S_1, \dots, S_T) \parallel P(Z_t)\right) = H \log 2 + \sum_{z=1}^{2^H} q(Z_t = z; S) \log q(Z_t = z; S).$$

The KL divergence can be expressed per token

$$\mathbb{D}_{\text{KL}}\left(q(Z_t; S_1, \dots, S_T) \parallel P(Z_t)\right) = H \log 2 + \sum_{z=1}^{2^H} q(Z_t = z; S) \log q(Z_t = z; S).$$

We use the free-bits method (Kingma et al., 2016), also per token

$$\mathcal{L}_{\text{KL}} = \frac{1}{T} \sum_{t=1}^T \max\left(0, \mathbb{D}_{\text{KL}}\left(q(Z_t; S_1, \dots, S_T) \parallel P(Z_t)\right) - \kappa\right).$$

# Outline

Understanding Transformers as Structure Induction Models

Challenges: Abstraction

Challenges: Compositionality

Challenges: Causal Latent Variables

Summary of Future Challenges

# Summary of Future Challenges

The structure of hidden representations will continue to be central to effective deep learning architectures for NLP and AI.

- ▶ Recent progress is rooted in unsupervised learning of linguistically appropriate structured representations.
  - ▶ Can we learn nodes of the structure (entities)?
  - ▶ Can we learn levels of representation?
- ▶ Can LLMs learn representations which capture fundamental properties of language?
  - ▶ Systematicity, yes.
  - ▶ Compositionality?
  - ▶ Entailment?
  - ▶ Multi-step reasoning?
- ▶ What architectural changes are needed to improve the inductive bias of LLMs?
  - ▶ Causal random choices within the latent representations?
  - ▶ ...