

Benchmarking LLMs

(for Particle Physics Understanding)

Student: Eugene Shalugin

Supervisor: Sascha Caron

Radboud University

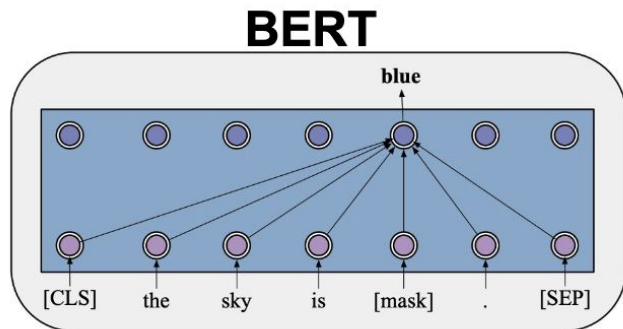
What is Understanding ?

The degree to which agent A scientifically understands phenomenon P can be determined by assessing the extent to which

- (i) A has a sufficiently complete representation of P;
- (ii) A can generate internally consistent and empirically adequate explanations of P;
- (iii) A can establish a broad range of relevant, correct counterfactual inferences regarding P.

(i-iii) can be measured, given a certain context (series of prompts) via what-, why, and w-questions respectively.

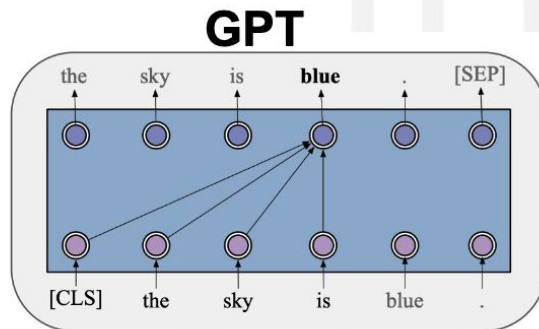
What is LLM ?



- Good at natural language understanding
- Maximizes the likelihood:

$$\mathcal{L}(\mathcal{X}) = \sum_{i=1}^m \log P([\text{Mask}]_i = y_i | \tilde{\mathcal{X}}; \Theta)$$

- P is modeled by transformer *encoder*
- $\tilde{\mathcal{X}}$: text after masking some tokens

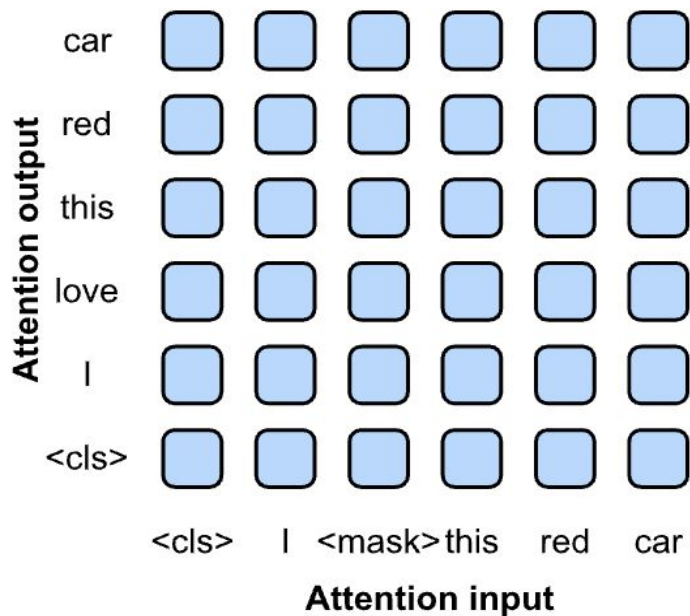
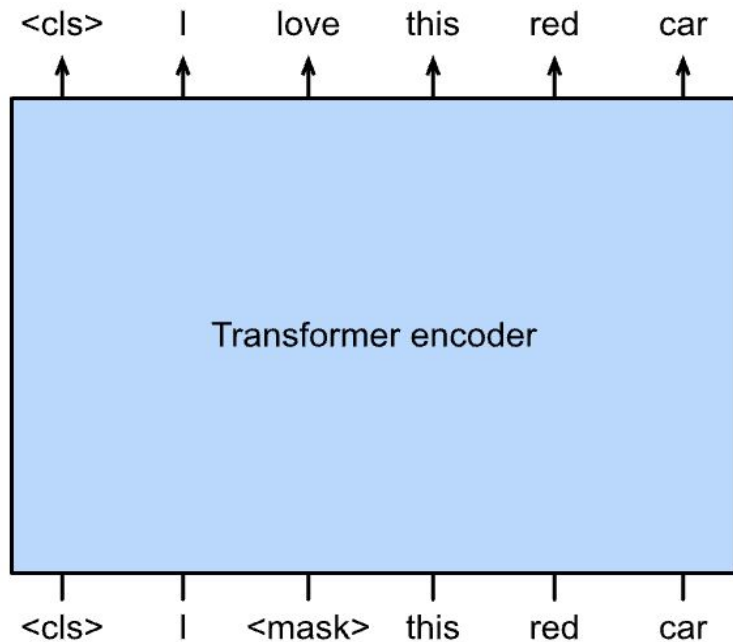


- Good at natural language generation
- Maximizes the likelihood:

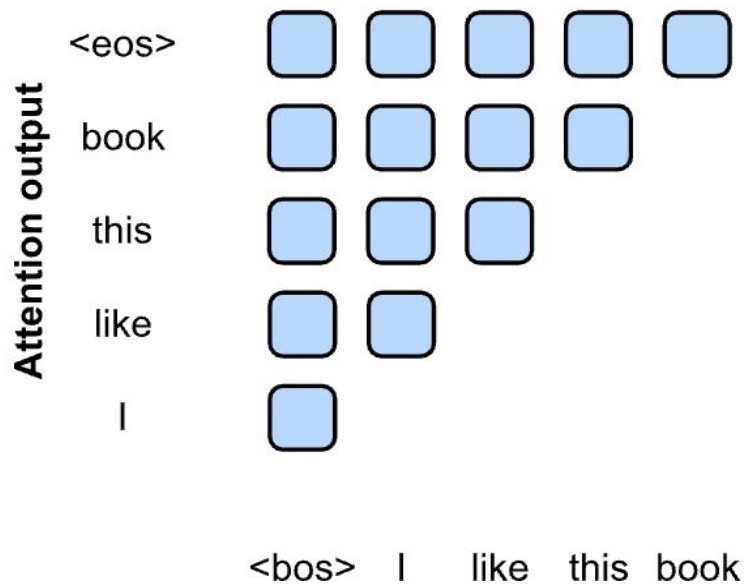
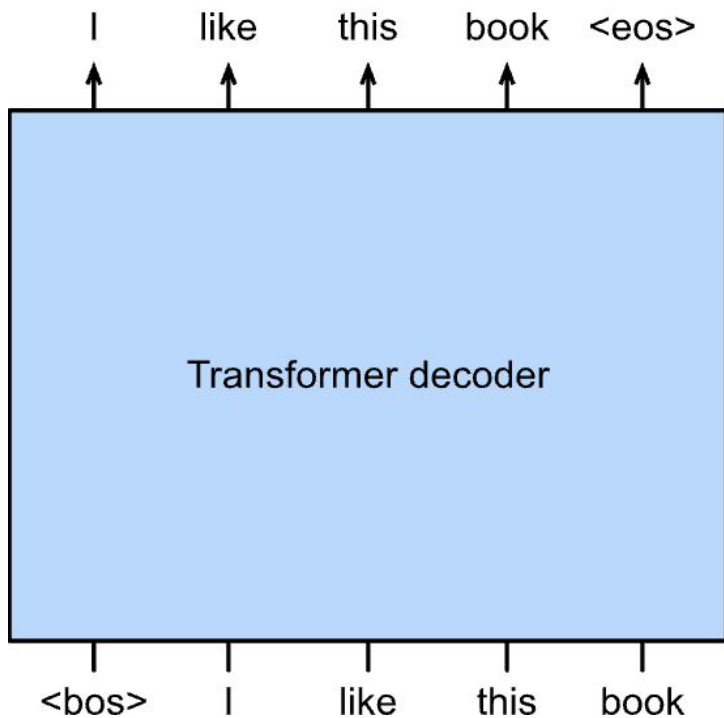
$$\mathcal{L}(\mathcal{X}) = \sum_{i=1}^{n+1} \log P(x_i | x_{i-k}, \dots, x_{i-1}; \Theta)$$

- P is modeled by transformer *decoder*
- *Self attention: left to right*

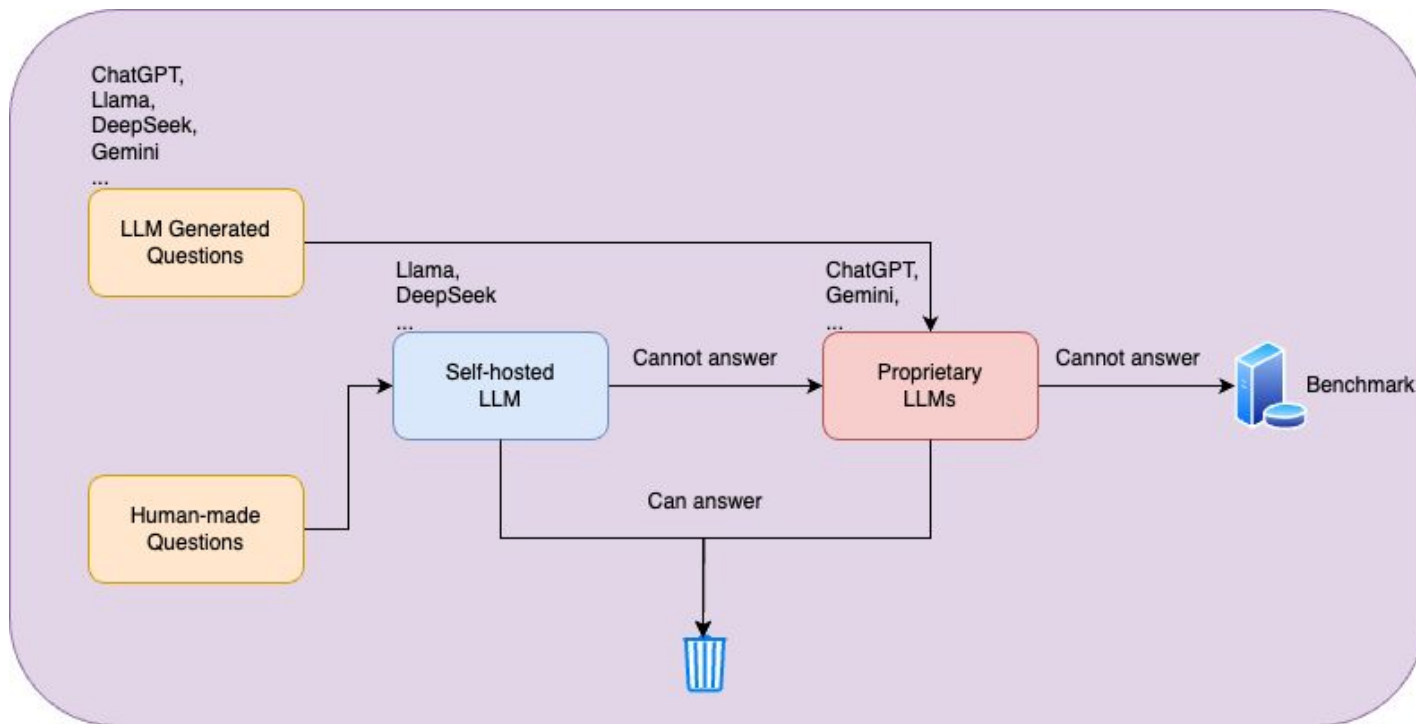
What is LLM – BERT



What is LLM – GPT



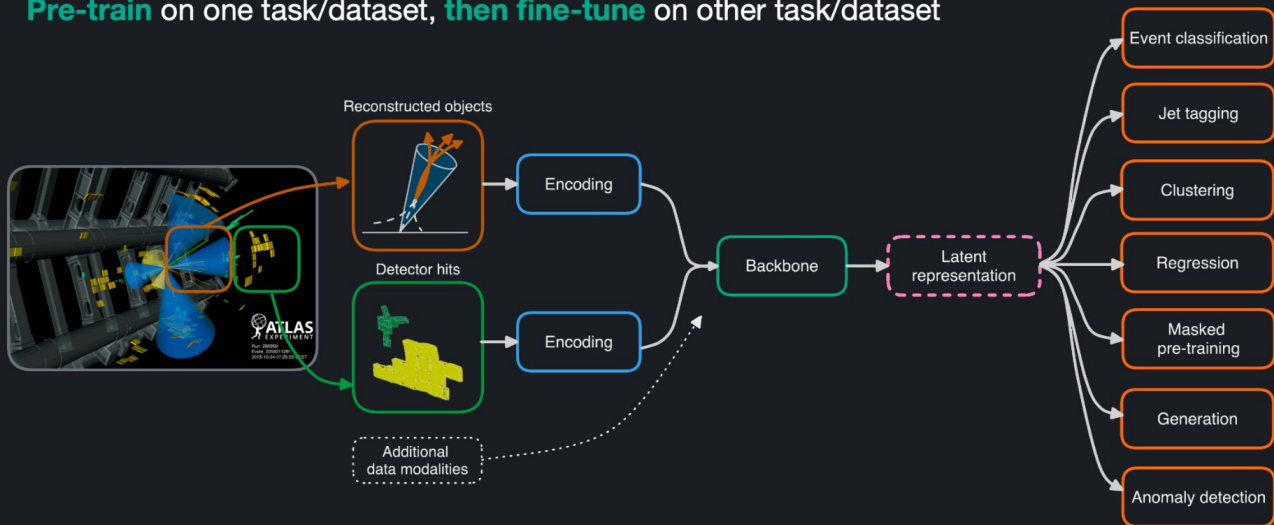
So what's with Benchmarking



Foundation Models

Foundation models for HEP

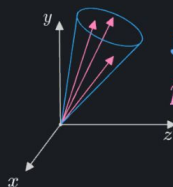
Pre-train on one task/dataset, then fine-tune on other task/dataset



Foundation Models

Our approach

Jet constituents with **continuous features**



$$\text{Jet} = \{\vec{p}_1, \vec{p}_2, \dots, \vec{p}_n\}$$
$$\vec{p}_i = (p_T, \eta^{\text{rel}}, \phi^{\text{rel}})$$

Constituents are **tokenized with a VQ-VAE**
(using the approach presented by Sam Klein earlier)

$$\text{Jet} = \{\text{start-token}, \text{token}_1, \dots, \text{token}_n, \text{end-token}\}$$
$$\text{token}_i = \text{integer value} \in [1, \dots, 8192]$$

Unsupervised pre-training of transformer backbone
on generative task (next-token prediction)

Fine-tuning to classification task:
Swap model head and copy over the
weights from the pre-trained backbone



Large (Language) Models for Physics

AI-driven Science

Event Analysis

Information Retrieval

Creative Inference

Hypothesis Generation

Cross-Domain Knowledge
Transfer

Supporting Anomaly
Detection

Automatic Data Annotation

Querying Experimental Data

Enhancing Experimental
Design

Contact



eugene.shalugin@ru.nl

sascha.caron@ru.nl