

Arrows of Time for Large Language Models

Vassilis Papadopoulos* Jérémie Wenger† Clément Hongler*

*EPFL †Goldsmiths, University of London

Why Arrows of Time? \rightsquigarrow Mathematical Models

Simple Model: Prime Multiplication

Language $\{3 \times 13 = 57 \quad 11 \times 13 = 143 \quad 11 \times 17 = 187 \dots p \times q = pq \dots\}$

Easier to learn to multiply than to factor $\rightsquigarrow \mathcal{D}_{CE}^{\rightarrow} > 0$

How can $\mathcal{D}_{CE}^{\leftrightarrow}$ become spontaneously > 0 ?

Learning and Sparsity: Linear Languages

$\mathcal{P}_A: \underbrace{[x_1, x_2, \dots, x_m]}_{m \text{ iid bits}}; \underbrace{[y_1, y_2, \dots, y_m]}_{m \text{ iid bits}} \in \mathbb{F}_2^m$

$y = A^{\rightarrow} x \quad x = A^{\leftarrow} y \quad (A^{\rightarrow})^{-1} = A^{\leftarrow}$

\mathcal{M}^{\leftarrow} learns \mathcal{P}_A easily starting from \mathcal{P}_B

$\leftrightarrow A^{\rightarrow} - B^{\rightarrow}$ sparse

Sparsity Symmetry Breaking

$A^{\rightarrow} - B^{\rightarrow}$ sparse \Rightarrow Typically $A^{\leftarrow} - B^{\leftarrow}$ less sparse

Communication Model: Alice, Bob, and Carol.

Suppose Alice, Bob, and Carol share a common language \mathcal{P}_B , with Alice and Bob using forward models and Carol using a backward model.

Now if Alice manages to learn \mathcal{P}_A easily from \mathcal{P}_B , $A^{\leftarrow} - B^{\leftarrow}$ should be sparse, so if she sends \mathcal{P}_A samples

Bob will learn \mathcal{P}_A easily; for Carol it will be harder. Language is "selected" to be easy forward, so backward is harder.

Next-Token Prediction:

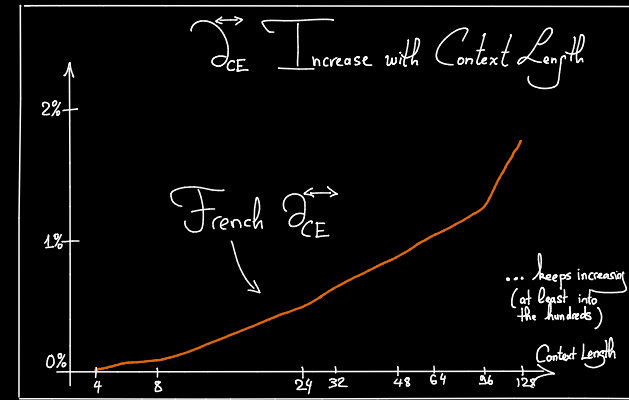
Once upon a time, there was a ? ...

\rightarrow Estimate $\mathbb{P}\{X_k = x_k \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}\} \forall k=1, \dots, n$

Previous-Token Prediction:

... ? and they lived happily ever after.

\rightarrow Estimate $\mathbb{P}\{X_k = x_k \mid X_{k+1} = x_{k+1}, \dots, X_n = x_n\}$



Arrow of Time: $\mathcal{D}_{CE}^{\leftrightarrow} = \frac{E[l_{CE}^{\leftarrow}] - E[l_{CE}^{\rightarrow}]}{\frac{1}{2}(E[l_{CE}^{\leftarrow}] + E[l_{CE}^{\rightarrow}])}$
measures the asymmetry in the estimability of languages by LLMs

Key Takeaways

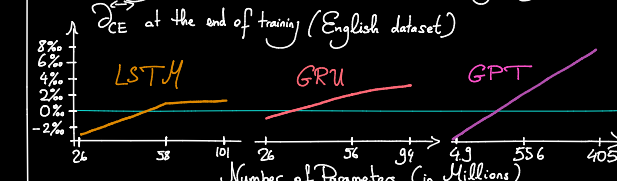
Arrows of Time are universal across languages, model architectures, model sizes, and training times, provided the datasets are large enough, and the models have enough parameters and training time.

While the effect is impressively robust, $\mathcal{D}_{CE}^{\leftrightarrow}$ is never very large (at most a few percent).

The effect size increases with the context length.

Artifacts due to the tokenization can be ruled out.

Universality across Architectures: $\mathcal{D}_{CE}^{\leftrightarrow} > 0$ as soon as models get large enough



Forward Model $\mathcal{M}^{\rightarrow}$:

$\mathbb{P}^{\rightarrow}\{X_1 = x_1, \dots, X_n = x_n\} = \prod_{k=1}^n \mathbb{P}^{\rightarrow}\{X_k = x_k \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}\}$

Backward Model \mathcal{M}^{\leftarrow} :

$\mathbb{P}^{\leftarrow}\{X_1 = x_1, \dots, X_n = x_n\} = \prod_{k=1}^n \mathbb{P}^{\leftarrow}\{X_k = x_k \mid X_{k+1} = x_{k+1}, \dots, X_n = x_n\}$

Two \mathbb{P} Factorizations: $\mathbb{P}\{X_1 = x_1\} \mathbb{P}\{X_2 = x_2 \mid X_1 = x_1\} \dots \mathbb{P}\{X_n = x_n \mid X_1 = x_1, \dots, X_{n-1} = x_{n-1}\} = \mathbb{P}\{X_1 = x_1\} \mathbb{P}\{X_2 = x_2 \mid X_2 = x_2\} \dots \mathbb{P}\{X_n = x_n \mid X_n = x_n\}$

Train a FW & BW copies of the same LLM (BW \Leftrightarrow FW on time-reversed dataset) \rightsquigarrow Do we have $\mathbb{P}^{\rightarrow} = \mathbb{P}^{\leftarrow}$?

Information-Theoretically: No Difference Between \mathbb{P}^{\rightarrow} and \mathbb{P}^{\leftarrow}

\hookrightarrow Do we see $l_{CE}^{\rightarrow} = l_{CE}^{\leftarrow}$?

Shannon's Experiments: Next- and Previous-Letter Prediction

Prediction and Entropy of Printed English
By C. E. SHANNON
(Manuscript received Sep. 21, 1948)

The idea of measuring the sum of cross-entropies in natural languages was pioneered by Shannon. He noted that this could also be done backwards...

Experiments were performed on human subjects; Shannon noted that to his surprise, they would perform worse predicting backwards but only slightly so.

Training: Minimize Cross-Entropy Losses

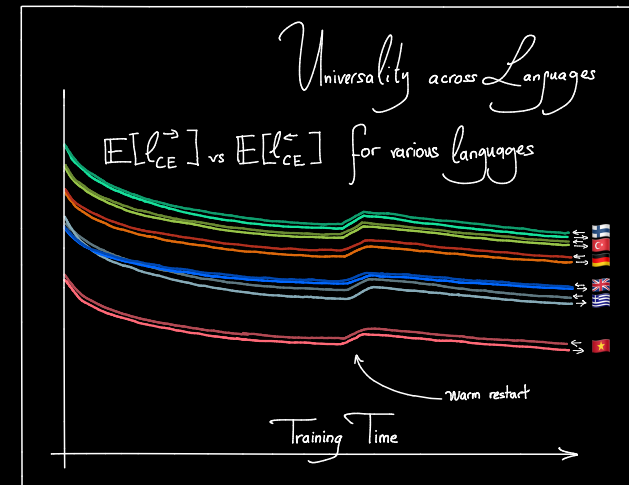
$l_{CE}^{\rightarrow} = \sum_{k=1}^n -\log \mathbb{P}^{\rightarrow}\{X_k = x_k \mid X_1 = x_1, \dots, X_{k-1} = x_{k-1}\}$

$= -\log \mathbb{P}^{\rightarrow}\{X_1 = x_1, \dots, X_n = x_n\}$

$l_{CE}^{\leftarrow} = \sum_{k=1}^n -\log \mathbb{P}^{\leftarrow}\{X_k = x_k \mid X_{k+1} = x_{k+1}, \dots, X_n = x_n\}$

$= -\log \mathbb{P}^{\leftarrow}\{X_1 = x_1, \dots, X_n = x_n\}$

\hookrightarrow if $\mathbb{P}^{\rightarrow} = \mathbb{P}^{\leftarrow}$ then we should have $l_{CE}^{\rightarrow} = l_{CE}^{\leftarrow}$



Open Questions & Perspectives

Arrows of Time in Code? Continuous Setting?

Arrows of Time in DNA? Link with Thermodynamics?

Arrows of Time \leftrightarrow Life? Link with Causality?

Quadratic Languages and Complexity? Stock Market Prices?

Small-Data Arrows of Time? Scaling Laws?