# Language Models and Cross-Entropy: Arrows of Time and the Life of Games

Clément Hongler

works with

Andrew Emil
Vassilis Papadopoulos
Jérémie Wenger

## Predictions, Scoring Rules, and Information

Setting : a random observation $i \in \{1, ..., n\}$

Prediction: $\vec{\pi} \in \Delta_n \leftarrow$ Simplex $\{\vec{\pi} \in [0,1]^n : \sum \pi_i = 1\}$

Scoring rule: $r: \{1, ..., n\} \times \Delta_n \to \mathbb{R}$

Upon $i$ being observed, reward $\vec{\pi}$ with $r(i, \vec{\pi})$

Example (bad!): $r(i, \vec{\pi}) = \pi_i$     "True prob."

Expected predictor's reward: $E(\vec{\pi}) = \sum_{i=1}^{n} p_i \cdot r(i, \vec{\pi})$

Proper scoring rule: $E(\vec{\pi})$ is maximized when $\vec{\pi} = \vec{p}$

Examples: quadratic: $r(i, \vec{\pi}) = -(1-\pi_i)^2 - \sum_k \pi_k^2 = -\|\vec{\pi} - \vec{\delta_i}\|^2$

quartic: $4\pi_i^3 - 3\sum_k \pi_k^4$   cross-entropy (xent): $r(i, \vec{\pi}) = \log \pi_i$

Proper scoring rule classification (Savage, ...)

$r(i, \vec{\pi}) = G(\vec{\pi}) + \langle \vec{\delta_i} - \vec{\pi}, \nabla G(\vec{\pi}) \rangle$     (G convex)

Xent is special: it is the only local proper scoring rule

$r(i, \vec{\pi})$ only depends on $\pi_i$

Xent leads us to information theory ...

# Language Models, Xent, and Compression

Given $n$ successive tokens ($\approx$ words) $x_1, ..., x_n \in V$ in a text, an autoregressive LLM outputs a prediction $(\pi_x^{n+1})_{x \in V}$ for the next token $x_{n+1}$

## Cross-Entropy (Xent) loss :  (negative reward)

$$-\sum_{n=0}^{N} \log \left(\pi_{x_{n+1}}^{n+1}\right)$$

Estimate of $\mathbb{P}\{X_{n+1} = x_{n+1} \mid X_1 = x_1, ..., X_n = x_n\}$

Estimate of $\log \mathbb{P}\{X_1 = x_1, ..., X_N = x_N\}$

The xent loss corresponds to the compressed size ($\approx$ shortest description) of a long text, if we use the LLM measure as an a priori: knowing the first $n$ tokens, we can use e.g. arithmetic encoding to encode the $n+1$-st, using (in expectation) $-\sum_{x \in V} P_x^{n+1} \log_2 \pi_x^{n+1}$ bits

Minimized when $\pi^{n+1} = P^{n+1}$

The pre-training process for LLMs consists in minimizing the xent over large sets of texts
The idea that the most concise description leads to intelligence is suggested by Occam's Razor, and the works of Shannon, Kolmogorov, Solomonoff, Chaitin
In practice, when the optimization is performed on transformers, the results are incredibly good!

Next - Token Prediction:
Once upon a time, there was a [ ? ] ...
$\rightarrow$ Estimate $\mathbb{P}\{X_k = x_k \mid X_1 = x_1, ..., X_{n-1} = x_{k-1}\} \, \forall k = 1, ..., n$

Forward Model $\mathcal{M}^{\rightarrow}$:
$$\mathbb{P}^{\rightarrow}\{X_1 = x_1, ..., X_n = x_n\} = \prod_{k=1}^{n} \mathbb{P}^{\rightarrow}\{X_k = x_k \mid X_1 = x_1, ..., X_{k-1} = x_{k-1}\}$$

Previous - Token Prediction:
... [ ? ] and they lived happily ever after.
$\rightarrow$ Estimate $\mathbb{P}\{X_k = x_k \mid X_{k+1} = x_{k+1}, ..., X_n = x_n\}$

Backward Model $\mathcal{M}^{\leftarrow}$:
$$\mathbb{P}^{\leftarrow}\{X_1 = x_1, ..., X_n = x_n\} = \prod_{k=1}^{n} \mathbb{P}^{\leftarrow}\{X_k = x_k \mid X_{k+1} = x_{k+1}, ..., X_n = x_n\}$$

# Information and Time Reversibility

## Can an LLM learn to speak backwards?

Train on data with a reversed time direction!

Useful for many things (reverse prompting, etc)

Train FW and BW models: same, except time

Xent answer: optimal compression is time flip-invariant

Probability answer:

$$\mathbb{P}\{X_1=x_1\}\,\mathbb{P}\{X_2=x_2\mid X_1=x_1\} \cdots \mathbb{P}\{X_N=x_N\mid X_{<n}=x_{<n}\}$$
$$= \mathbb{P}\{X_N=x_N\}\,\mathbb{P}\{X_{N-1}=x_{N-1}\mid X_N=x_N\}\cdots \mathbb{P}\{X_1=x_1\mid X_{>1}=x_{>1}\}$$

Compare forward and backward LLMs with same data

- Information-theoretically: no difference
- If we memorize dataset: no difference
- Naively: we can't speak backwards, so that LLMs probably can't either
- Shannon: interesting question, unclear
- Google rumors (circa 2014): no difference
- Some (not well-cited) papers: backward easier

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

→ 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

$$\vec{\ell}_{CE} = \sum_{k=1}^{n} -\log \overrightarrow{\mathbb{P}}\{X_k=x_k \mid X_1=x_1, \ldots, X_{k-1}=x_{k-1}\}$$

$$= -\log \overrightarrow{\mathbb{P}}\{X_1=x_1, \ldots, X_n=x_n\}$$

$\log \overrightarrow{\mathbb{P}}:$ 

once upon a time, lived a wise old woman who knew the secrets of the forest

$\log \overleftarrow{\mathbb{P}}:$

$$\overleftarrow{\ell}_{CE} = \sum_{k=1}^{n} -\log \overleftarrow{\mathbb{P}}\{X_k=x_k \mid X_{k+1}=x_{k+1}, \ldots, X_n=x_n\}$$

$$= -\log \overleftarrow{\mathbb{P}}\{X_1=x_1, \ldots, X_n=x_n\}$$

↪ if $\overrightarrow{\mathbb{P}} = \overleftarrow{\mathbb{P}}$ then we should have $\vec{\ell}_{CE} = \overleftarrow{\ell}_{CE}$
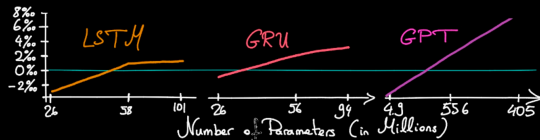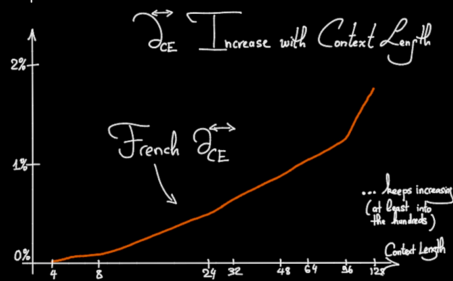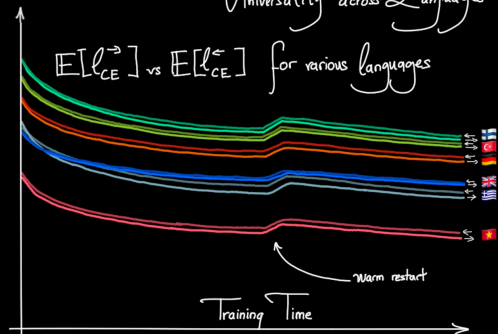
# Universal Arrows of Time for LLMs
(w/ V. Papadopoulos & J. Wenger)

Compare FW and BW xent losses at end of training:

- For all human languages $FW$ loss $<$ $BW$ loss

- The difference increases with the context length

- For small models pre-2017, the effect is tiny or reversed

- As the model sizes increase, the AoT becomes large
(Work by Zhong, Bai, Gu, Zhang, Gu, Abbé, Bengio, Jaitly)

- The effect is related to the semantics, not the syntax

- The effect strength depends on the language (why?)

- The same effect can be observed on code (w/ Y. Romaniv)

- This effect cannot be observed on DNA

- The effect is stronger on LLM-generated data


Where do AoTs come from?

## Universality across Languages

$E[\ell_{CE}^{\rightarrow}]$ vs $E[\ell_{CE}^{\leftarrow}]$ for various languages



Warm restart

Training Time

## $\partial_{CE}^{\leftrightarrow}$ Increase with Context Length



2%

1%

0%

French $\partial_{CE}^{\leftrightarrow}$

...keeps increasing
(at least into
the hundreds)

Context Length

4  8    24 32   48 64   96 128



8‰
6‰
4‰
2‰
0‰
-2‰

LSTM          GRU          GPT

26   58  101    26   56   94    49   556   405

Number of Parameters (in Millions)

# AoTs and Computational Hardness

## A simple-to-understand example:

Dataset: a multiplication table $p \times q = pq$ for primes $p < q$

- The FW model predicts RHS from LHS by multiplying
- The BW model predicts LHS from RHS by factoring

Modern LLMs can learn to multiply well; factoring is hard

Numerical Example $\quad p < q, \quad p \cdot q < 10^5$

|    | P    | q    | pq (reversed) |
|----|------|------|---------------|
| FW | 8.98 | 8.67 | 4.55          |
| BW | 0.02 | 8.41 | 21.56         |

Similarly $pq = p \times q$ exhibits a reverse direction...

## Why do FW AoTs arise spontaneously?

# AoTs via Sparsity Symmetry Breaking

Idea: the reverse of a sparse circuit is (generically) sparse, but not as sparse (hard to prove, but intuitive)

Linear language dataset: "$x \longleftrightarrow y$" where $x, y \in \mathbb{F}_2^m$ are made of $m$ i.i.d. bits $x = f^{\leftarrow} y, \; y = f^{\rightarrow} x$, for fixed bijective matrices $f^{\rightarrow}, f^{\leftarrow} : \mathbb{F}_2^m \rightarrow \mathbb{F}_2^m$

- The LHS of "$\longleftrightarrow$" determines the RHS and vice versa
- If a random $f^{\rightarrow}$ is sparse $f^{\leftarrow}$ is typically less sparse

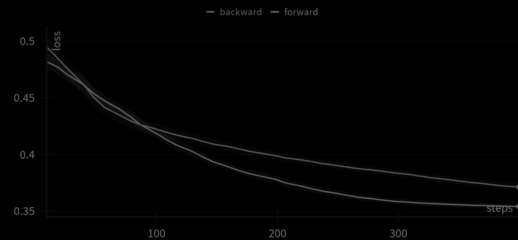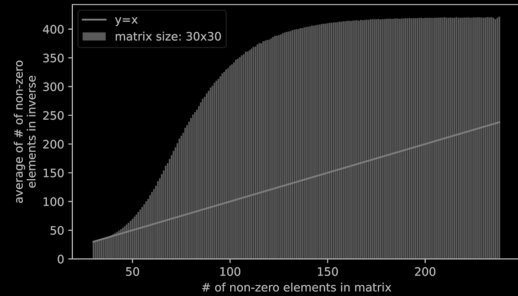Now: sparsity $\longleftrightarrow$ ease of learning
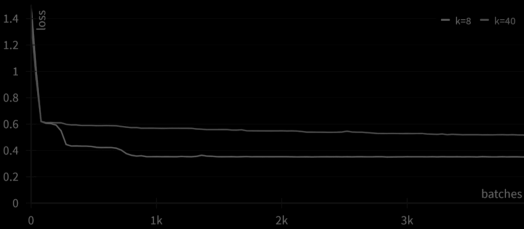
So: if we condition a linear language to be sparse FW, it is less sparse BW

Similarly, if we perform an update to a lin. lang. that is FW sparse, it is less BW sparse

Suppose Alice & Bob are FW agents with a common language $L$, and Carol is a BW agent who knows $L$. Now if Alice updates $L$ sparsely, the update for Carol will typically be less sparse $\rightsquigarrow$ AoT!

## Charts (left)

**Top chart:** loss vs batches — legend: k=8, k=40. Y-axis: 0 to 1.4. X-axis: 0, 1k, 2k, 3k.

**Middle chart:** average of # of non-zero elements in inverse vs # of non-zero elements in matrix. Legend: y=x, matrix size: 30x30. Y-axis: 0 to 400. X-axis: 50, 100, 150, 200.

**Bottom chart:** loss vs steps. Legend: backward, forward. Y-axis: 0.35 to 0.5. X-axis: 100, 200, 300.

## AoT Mysteries and Perspectives

- Can we make sense of the following idea:

Theoretically:

$$\log \mathbb{P}\{\text{initial configuration}\} + \sum_{\text{step}} \log \mathbb{P}\{\text{step} \mid \text{past}\}$$

$$= \log \mathbb{P}\{\text{final configuration}\} + \sum_{\text{step}} \log \mathbb{P}\{\text{step} \mid \text{future}\}$$

(Covered by entropy creation?)     (Covered by algorithmic AoT?)

- Can we find algorithmic AoTs with small datasets?
- Do we have AoTs in animal communication?

Maybe not?

- Can algorithmic AoTs arise in physical settings?
- Is the manifestation of AoTs a sign of "life"?
- Are there systems exhibiting reverse AoTs?
- Can we unify with the vision with diffusion models?