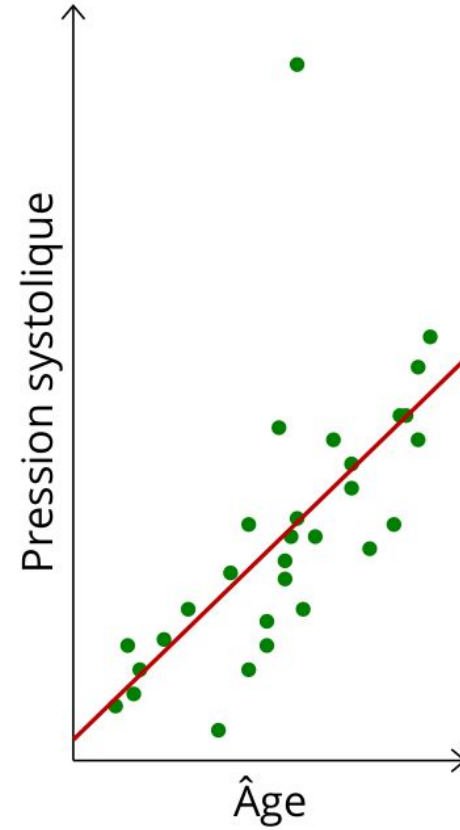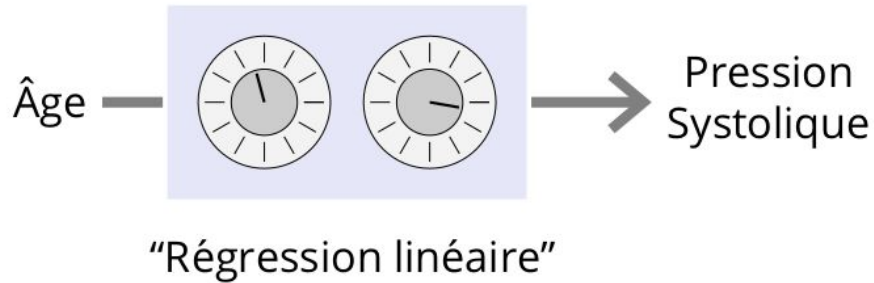# From deep learning to AI

NRB AI & Data Xperience, October 15

Prof. Damien Ernst, <u>Prof. Gilles Louppe</u>, Lize Pirenne
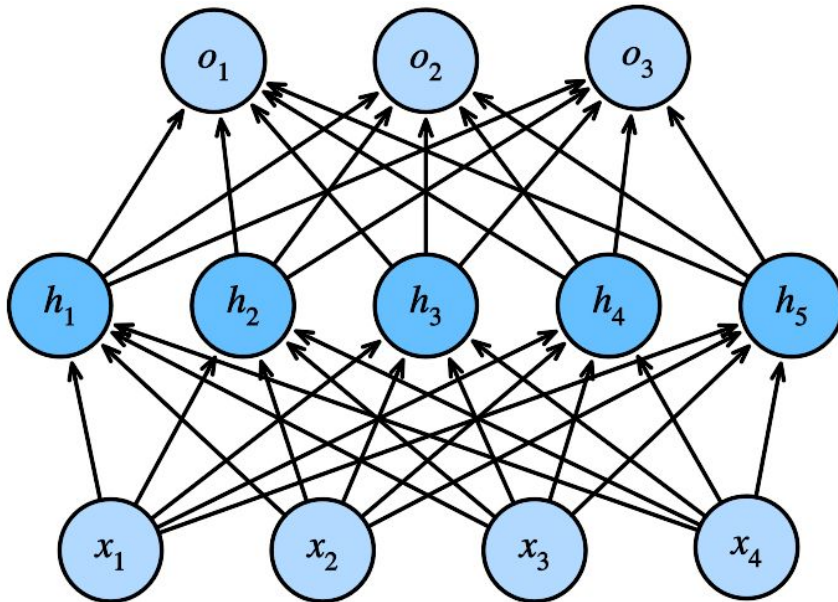
LIÈGE université

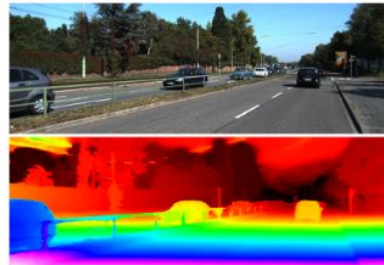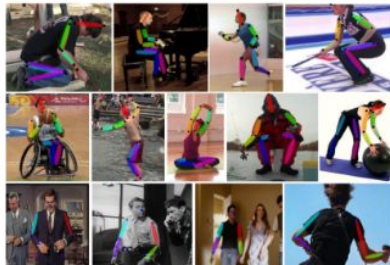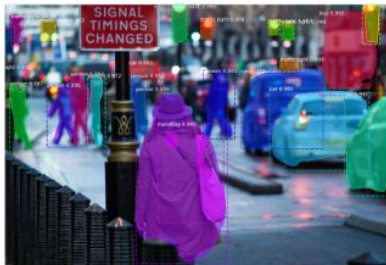Âge ── "Régression linéaire" ──▶ Pression Systolique

The **machine learning** approach to problem-solving.

Deep learning scales up the machine learning approach by

- using larger models known as **neural networks**,
- training on **larger datasets**,
- using **more compute** resources.

Specialized neural networks can be trained to achieve **super-human performance** on many complex tasks that were previously thought to be out of reach for machines.



(Top) Scene understanding, pose estimation, geometric reasoning.
(Bottom) Planning, Image captioning, Question answering.

Neural networks form **primitives** that can be transferred to many domains.



(Top) Analysis of histological slides, denoising of MRI images, nevus detection.
(Bottom) Whole-body hemodynamics reconstruction from PPG signals.

# The breakthrough

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

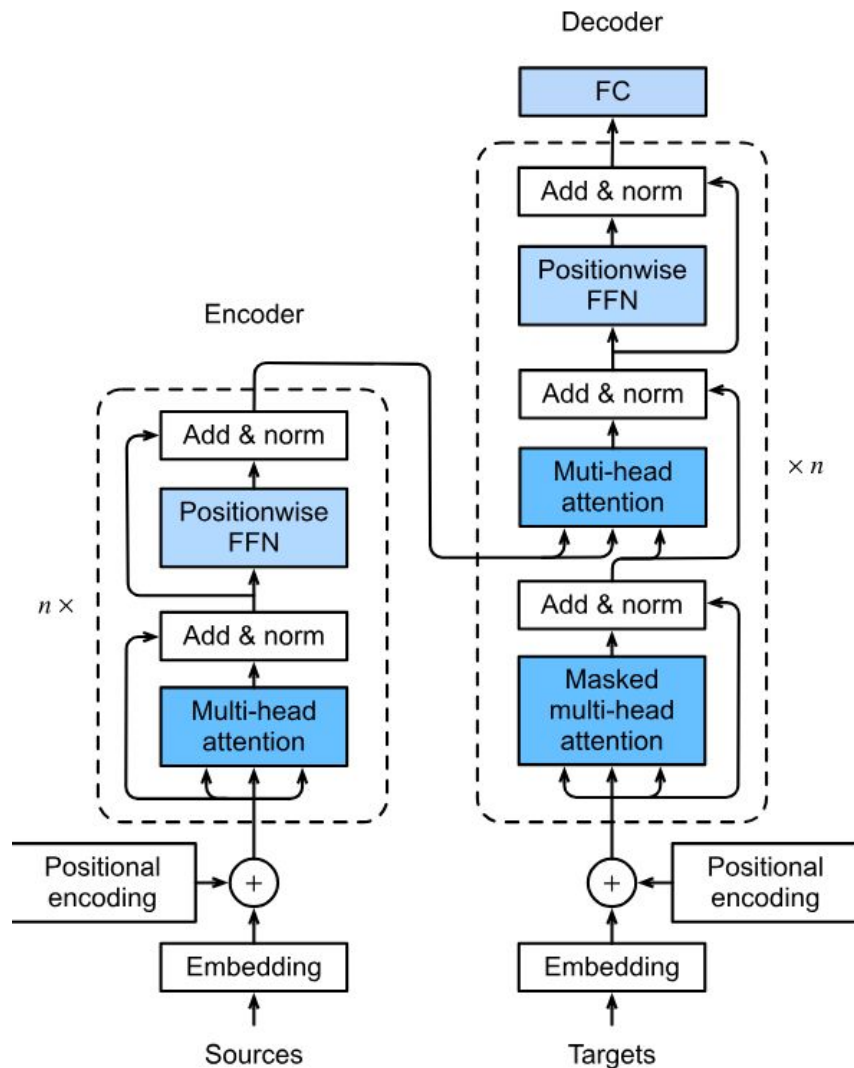Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu
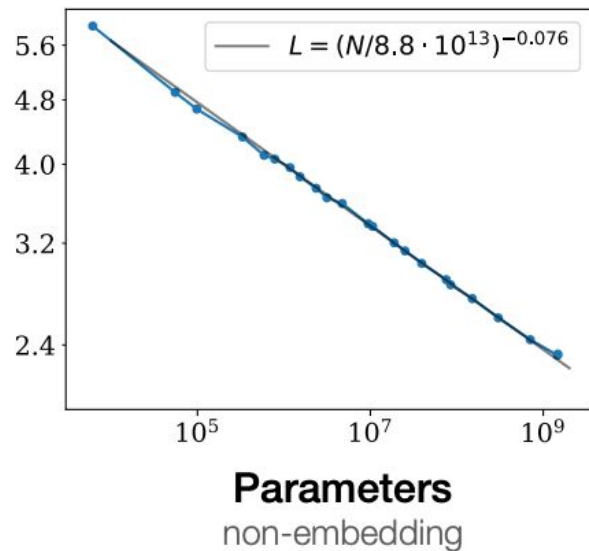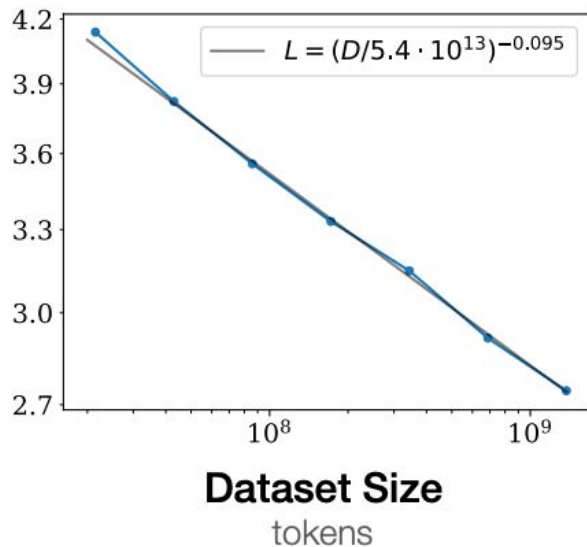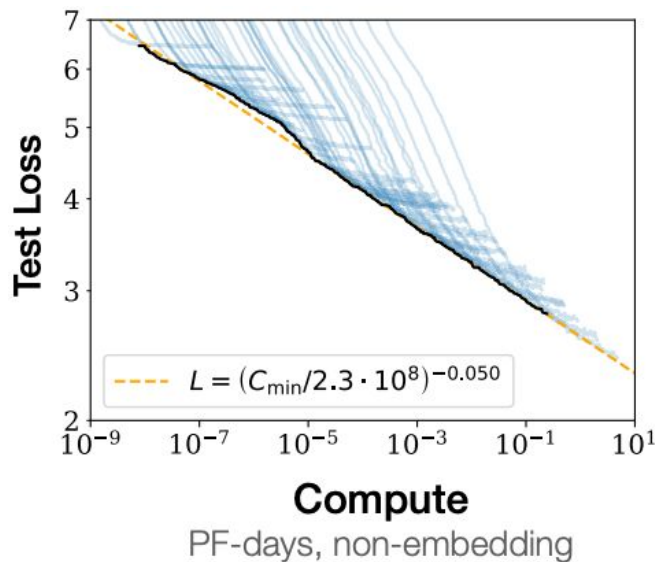
Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Vaswani et al, 2017.

Decoder

FC

Add & norm

Positionwise FFN

Add & norm

Muti-head attention

Add & norm

Masked multi-head attention

Encoder

Add & norm

Positionwise FFN

Add & norm

Multi-head attention

$n \times$

$\times n$

Positional encoding

Embedding

Sources

Positional encoding

Embedding

Targets

A **brutal simplicity**:

- The more data, the better the model.
- The more parameters, the better the model.
- The more compute, the better the model.

Conversational assistants (Anthropic, 2024)

Code assistants (Cursor, 2024)

Multi-modal assistants (OpenAI, 2024)

**New opportunities**

If Bob knew all of human knowledge and understood* everything you do, say or write…

- What questions would you ask Bob?
- What would do with Bob if you could hire him in your company?
- What if all your employees could talk with Bob?
- What if all your customers could talk with Bob?
- What if Bob was also working for all your competitors?

# Limitations

- LLMs give the impression of understanding.
- LLMs show no evidence of formal reasoning.
- LLMs are fragile. Rephrasing questions alter results significantly.
- LLMs do not know what do they not know.
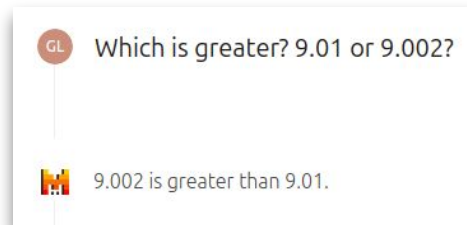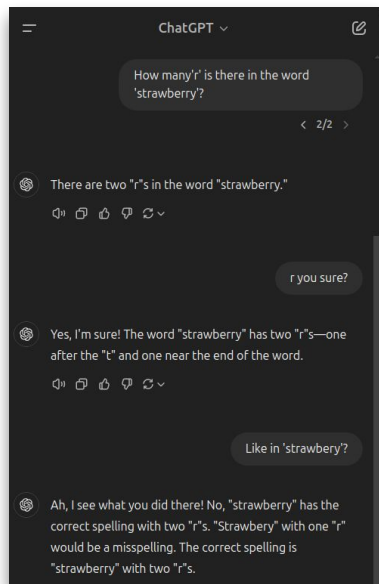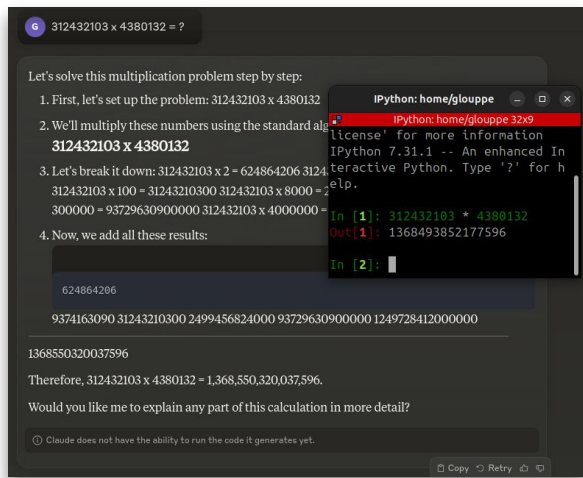


**Yann LeCun** ✔ ∞
@ylecun

Worth repeating:
Do not confuse retrieval with reasoning.
Do not confuse rote learning with understanding.
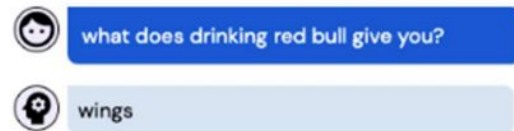Do not confuse accumulated knowledge with intelligence.

LLMs can be made more robust and helpful if the **rationale** of their answers can be **extracted or explained**.

what does drinking red bull give you?

wings

Page: Red Bull

Red Bull's slogan is "it gives you wings". The product is strongly marketed through advertising, tournament sponsorship, sports team ownerships, celebrity endorsments and with its record label.
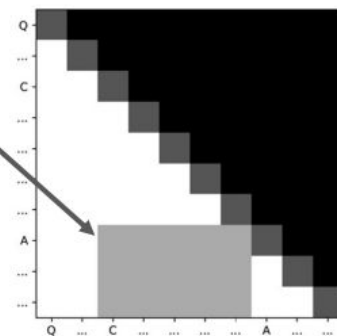
Our solution was to use the inner workings of an LLM to see what it pays attention to. More specifically, we were analysing the attention matrix to find the sentences in the context needed for generating the answer.

Pirenne et al, 2024.

~ Auto-correlation of the input

Rationale In black

Irrelevant in grey

Layer 8, Head 6

### Question: What is Sauvignon blanc?
### Context: Sauvignon blanc is a green-skinned grape variety that originates from the city of Bordeaux in France. The grape most likely gets its name from the French words sauvage ("wild") blanc ("white") due to its early origins as an indigenous grape in South West France. It is possibly a descendant of Savagnin. Sauvignon blanc is planted in many of the world's wine regions, producing a crisp, dry, and refreshing white varietal wine. The grape is also a component of the famous dessert wines from Sauternes and Barsac. Sauvignon blanc is widely cultivated in France, Chile, Romania, Canada, Australia, New Zealand, South Africa, Bulgaria, the states of Oregon, Washington, and California in the US. Some New World Sauvignon blancs, particularly from California, may also be called "Fumé Blanc", a marketing term coined by Robert Mondavi in reference to Pouilly-Fumé.

The research challenges of deploying LLMs for sensitive applications (e.g., medical, legal, financial) highlight **opportunities for mutually beneficial collaborations between the industry and academic**. (Many others exist, beyond LLMs!)

Benefits **for academia**:

- Access to relevant real-world problems.
- Access to industry-grade computing resources for research.
- Access to real-world data.
- Financial support.



Benefits **for industry**:

- Solutions for problems requiring extensive research efforts.
- Direct access to up-to-date knowledge in AI.
- Opportunities to work with brilliant minds and potential future employees.

# NRB - ULiège Research Chair

- 4-year funding for research on **AI for the industry**.
- Prof. Damien Ernst (LLMs, Reinforcement Learning, Energy)
  Prof. Gilles Louppe (Deep Learning, Generative AI, Digital Twins, AI4Science)
- Joint teams, regular meetings and shared objectives to encourage synergies.



(Left) Prof. Louppe and Prof. Ernst, (Right) Laurence Mathieu, CEO of NRB.