# Ia Artificial Intelligence all we need?

Daniele Bonacorsi (INFN e Università di Bologna)

INFN BOLOGNA
Istituto Nazionale di Fisica Nucleare
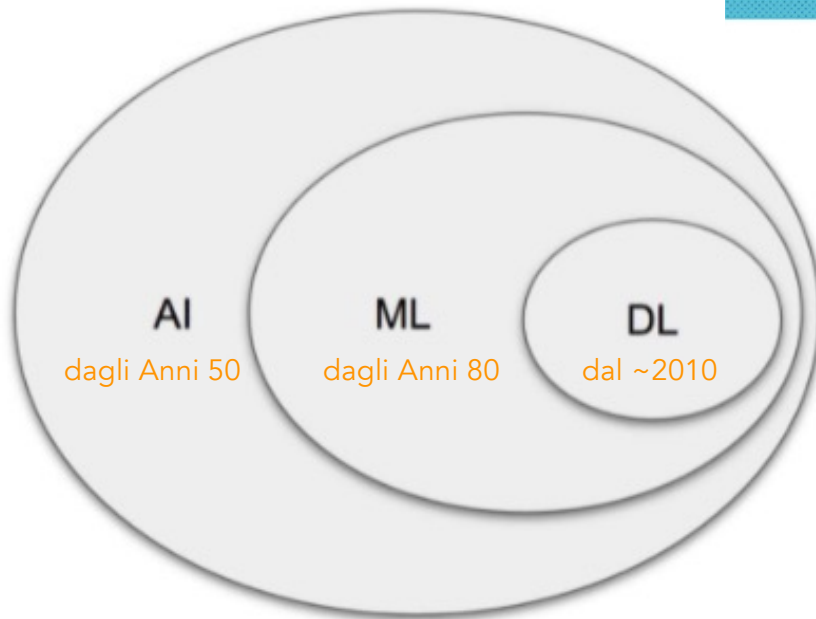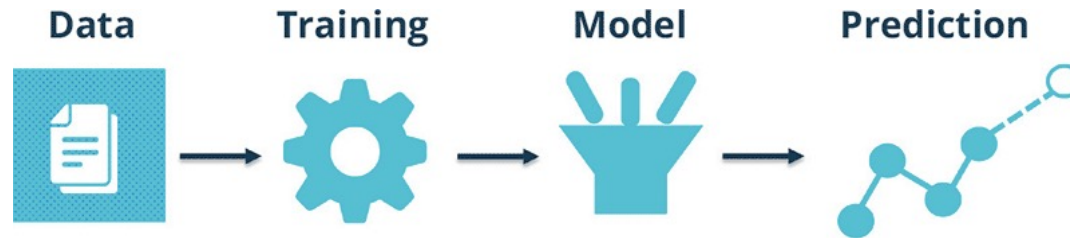Sezione di Bologna

INSIEME NELLE GFIDE PER LA SCIENZA DI DOMANI

LECCE
13e14 GIUGNO 2024
CINEMA MULTISALA MASSIMO

Istituto Nazionale di Fisica Nucleare
Piano Triennale
2025 | 2027

# AI vs ML vs DL

AI = **Artificial Intelligence**
ML = **Machine Learning**
NN = **Neural Network**
DL = **Deep Learning**



**Data** → **Training** → **Model** → **Prediction**

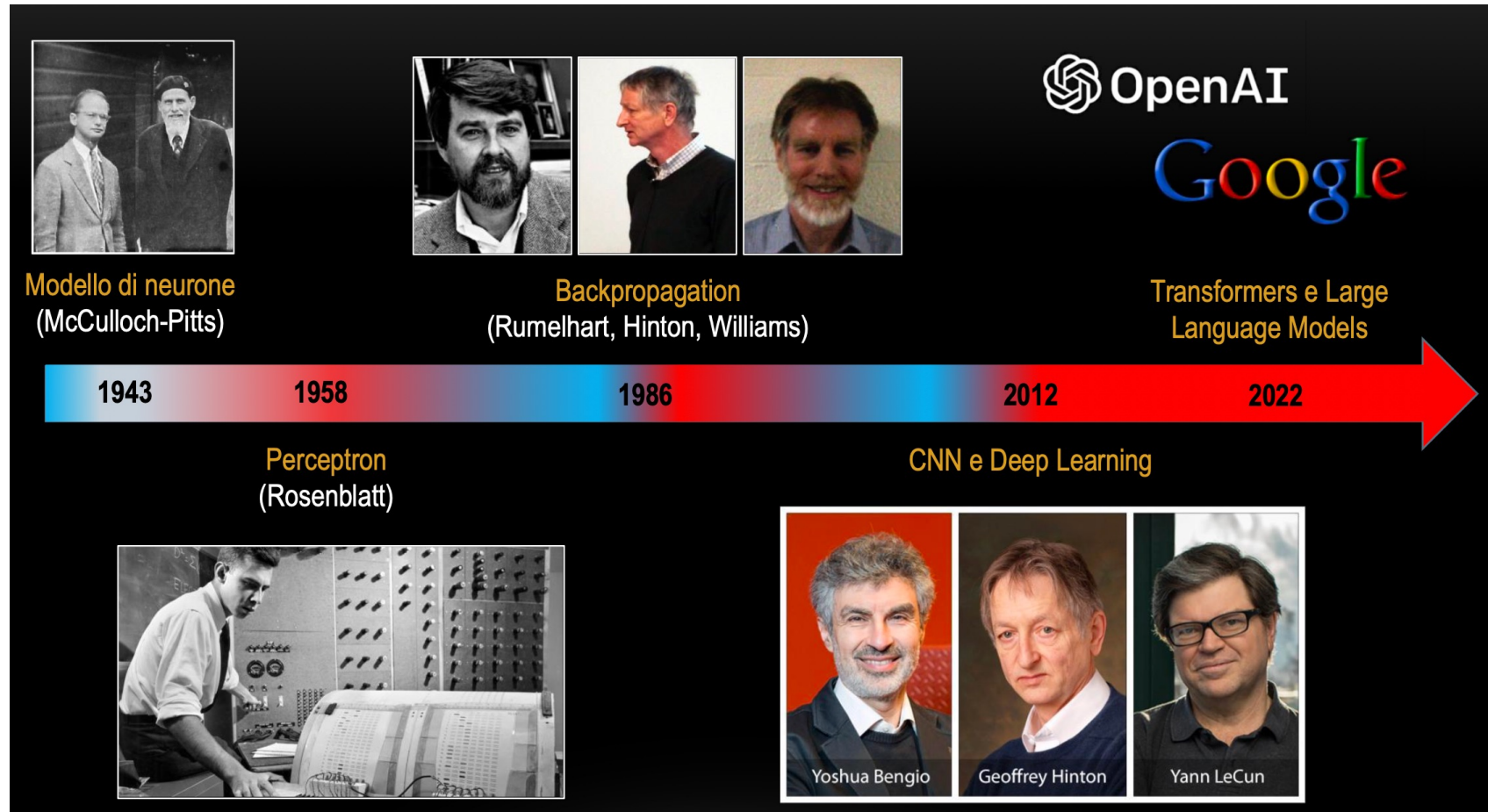AI — dagli Anni 50
ML — dagli Anni 80
DL — dal ~2010

<u>Objective</u>: extract "actionable insight" from (big) data

Choose an **algorithm**, perform its "**training**" on data ("attributes" vs "**features**") to extract "**parameters**" with optimisation techniques (e.g. "**gradient descent**") that minimise the errors of the model on the observations ("**cost function**"), in a process governed by "**hyper-parameters**" tuning
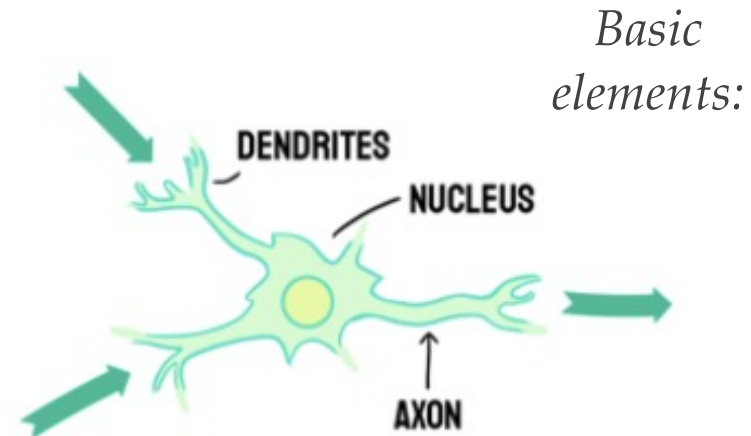
<u>Result</u>: a **ML model** to be applied to previously unseen data
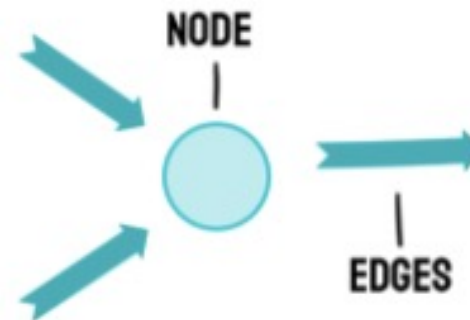
→ "**data-driven modelling**"

# Neural Networks: milestones



Credits: D. Maltoni

D. Bonacorsi

# Neural Networks

*Basic elements:*

Biological Neural Networks

**Artificial** Neural Networks

D. Bonacorsi

# A brief history of NNs
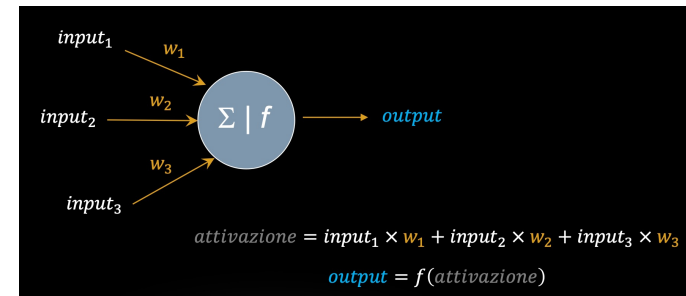
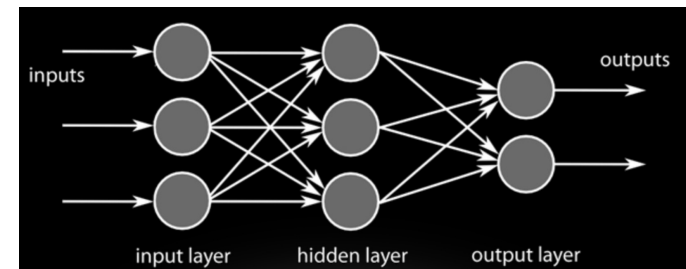First models of **artificial neuron** (McCulloch, Pitts, 1943)

- Neuron behaviour (i.e. info passing) are determined by "weights": initially randomly set, they are modified during "training"

Towards **neural networks**: "Perceptron" (Rosenblatt, 1958)

Multi Layer Perceptron (MLP)

- Groups of neurons organized in **layers** (**input** layer, **output** layer, and one or more intermediate "**hidden**" layers, each with one or more neurons)

**Back-propagation** algo (Rumelhart, Hinton, Williams, 1986)

- an effective technique to train them (i.e. setting weights' values)

From shallow NN to Deep NN → **Deep Learning**



$attivazione = input_1 \times w_1 + input_2 \times w_2 + input_3 \times w_3$

$output = f(attivazione)$

D. Bonacorsi

# E.g. "supervised learning": ML training vs inference

# AI's past winters and current spring

Past "**winters**" of AI

- '60: **shallow NN** hard to train

- '90: Support Vector Machines (**SVM**), Boosted Decision Trees (**BDT**), ..

- 2000+: advanced **deep NN** architectures

Current "**spring**" explainable by:

- **"Big data"**

- **Technology** + **ML research**

- (**cloud and accessibility**)

# ML-based publications in science



Extremely large adoption in incredibly short times, towards high level of **pervasiveness**

B. Blaiszik, "2021 AI/ML Publication Statistics and Charts", 10.5281/zenodo.7057437

ML in **data acquisition** and **trigger**

- Bkg and trigger rate reduction

- Signal specific trigger paths

- Anomaly detection in data taking

- Unsupervised new physics mining

E.g. LHC experiments' trigger is a strong "driver" for high-performances ML applications

- Next-gen trigger systems → real-time reconstruction → real time analysis

Challenge is the trade-off between **algorithmic complexity** and the **performances** achievable under severe time constraints in inference



40 MHz → L1 Trigger → 100 KHz → High-Level Trigger → 1 KHz 1 MB/event → Offline

*[ **CMS** - credits: M.Pierini ]*

# Pervasive ML in HEP [2/3]

ML in Event **Simulation**

The production of simulated events (full/fast simulation) is extremely intense from the computation standpoint (up to the point it might impact the physics reach of the experiments). ML can help to reduce such load

- Calorimeter shower surrogate simulator

- Analysis level simulator

- Pile-up overlay generator

- Monte Carlo integration

- ML-enabled fast-simulation

- Invertible full-simulation (probabilistic programming, …)

- …

ML in Event **Reconstruction**

Online/offline reconstruction might be partially replaced by surrogate models (approximate → faster) or by new algorithms (that might offer unprecedented performances)

- Charged particle tracking (GraphNN, vertexing, …)

- Calorimeter reconstruction (local, clustering, …)

- Particle flow (GraphNN, …)

- Particle identification (boosted jets, isolation, …)

- Pileup mitigation

- Energy regression (end-2-end, …)

- …

10

# Pervasive ML in HEP [3/3]

## ML in **Data Analysis**



Data: IML

**Legend:**
- ● W/Z tagging
- ● H→b bar
- ● quarks and gluons
- ● top quark tagging
- ● strange jets
- ● b-tagging
- ● flavor physics
- ● BSM particles and models
- ● particle identification
- ● neutrino detectors
- ● direct Dark Matter detectors
- ● cosmology, astroparticle, and cosmic ray physics
- ● tracking
- ● heavy ions / nuclear physics

## ML in **Computing Operations**

Application of ML to **non-collision (meta-)data** might help to increase efficiency and reduce the need omg personpower in Ops, e.g. automating specific tasks, developing intelligent/adaptive systems, ultimately acting on the full chain - from data collection to data analysis - and make it more agile

- Detector control

- Data quality monitoring

- Operational intelligence

- Predictive maintenance

- …

11

# ML/DL in HEP

"Traditional" ML



ML in HEP started by using domain knowledge to perform **feature extraction/engineering**

- HEP physicists design high-level features, and send them in punt to traditional ML "shall" algorithms

# Particle id, energy resolution, e oltre..

**BDT** usati per il learning delle energy corrections usando tutte le info disponibili nei vari sensori calorimetri

- es. **CMS**: energy sum, recupero bremsstrahlung con supercluster, inclusione della pre-shower energy, poi **energy correction con algo multivariato**





**LHCb** usa delle **NN** allenate su $O(30)$ features da tutti i sottosistemi, ciascuna allenata a identificare uno specifico tipo di particella

- → **mis-ID bkg/particle ridotta di ~3x (e oltre..)**

*[ 2015 ECAL detector performance plots, CMS-DP-2015-057. Copyright CERN, reused with permission ]*

# Uso di ML in scoperta e studio delle proprietà Higgs

Ruolo chiave del ML nella <u>scoperta del bosone di Higgs</u> prima delle attese

- es. diphoton analysis di **CMS**, **BDT** usati per migliorare la risoluzione e selezionare/categorizzare eventi

- → <u>**sensitività aumentata di un ammontare equivalente al ~50% di dati raccolti in più**</u>





*[courtesy M.Pierini]*

Studio delle <u>proprietà H</u>: es. $\tau$ leptons, **ATLAS** su 6 regioni cinematiche distinte, training di una BDT in ciascuna, con 12 features

→ <u>**sensitività migliorata del ~40% rispetto a un approccio non-ML**</u>

[1] JHEP 04 (2015) 117

# Test di alta precisione del MS

**CMS** e **LHCb**: evidenza per il decadimento raro $B^0_s \rightarrow \mu^+\mu^-$ con analisi combinata [1]

- es. **BDT** usati per ridurre la dimensionality dello spazio delle features, poi analisi spettro massa in BDT response bins → decay rate consistente con predizioni SM con precisione ~25% → constraints a estensioni SM

- → es. **per avere la stessa sensitività senza ML, LHCb come singolo esperimento avrebbe dovuto raccogliere ~4x più dati**



*Mass distribution of the selected $B^0 \rightarrow \mu^+\mu^-$ candidates with BDT > 0.5 [2]*

[1] Nature 522 68–72 (2015)
[2] *Phys.Rev.Lett.* 118 (2017) 19, 191801

# ML/DL in HEP



"Traditional" ML    Seeking DL solutions

Since some years, ML (DL) in HEP seeks for more advanced techniques, e.g. **deep NNs**

- **Use all the features space at its full dimensionality to train deep NN - no more manual feature engineering**

→ estract best from data, and do so by exploiting any architecture that might work for a given use-case (e.g. input from CV and NLP solutions..)

# Convolutional Neural Networks (**CNN**)

**CNNs** offer translational-invariant feature learning, robustness against noise, versatility in application to a variety of domains

- Extremely vast zoo of architectures! Primary target: computer vision

- are based on sequences of convolutional and pooling layers, and

**Industry**:
Large adoption in computer vision applications
(e.g., self-driving cars, ..)

**HEP**:
3D imaging in detectors, event classification,
automation of hist checking (e.g. data quality), ...



General tactics: (TPCs, CALOs..): represent your data as 2D/3D images (even 4D w/ timing info)

→ **problem casting into a computer vision task**

17

# "HEP is so different from other applications..". Are we?



Rivelazione di tracce di **neutrini** su cosmic background events (metodo: **CNN**)



Rivelazione di **aeroporti** da immagini satellitari (metodo: **CNN**)

# Recurrent Neural Networks (**RNN**)

**RNN** allow to handle variable-length inputs and process time-series, accumulating and using together info at various times in the sequence

- Based on "recurrent neurons" (backward-pointing connections)

- A variety of application in time-series of all kins, e.g. language translation, ..

**Industry**:
handling "time series"
(audio, video, natural language processing)

**HEP**:
Classifiers capable to process variable-length signals of different nature (tracks, particles in jets, etc) - ample application in astro-particle physics

# Autoencoder (**AE**)

**AE** is a "data-specific" compression algorithm, able to reduce dimensionality and extract "the juice" of an input

- a feed-forward (un/self-supervised) NN capable to encode the input into a reduced-dimensionally representation ("latent space") and decode it in output

**AEs in Industry**:
dimensionality reduction (like PCA),
clustering, denoising, …



**AEs in HEP**:
anomaly detection (intestìni events are those whose decoding in output is dtstant from the input, according to a given metric



Potenziale strumento per scoprire nuova fisica in modo "unsupervised"

22

# (V)AE per "new physics mining" - at LHC and beyond

A **Variational AE** has been proposed (CMS) for "new physics mining"

- Traing on known SM processes, build threshold to identity "anomalous" (i.e. interesting: BSM?) events

- Treat them as outliers, save them (no trigger kill!), build a catalog for further inspection

- Model-independence:  training not dependent from specific new physics signatures → assumptions-free

- Might be complementary to classical methods, i.e. model-dependent hypothesis testing

- topologie ricorrenti nel catalogo possono ispirare search focalizzate, e anche costruzione di modelli teorici

- Target: up to the trigger level..

Note: going from **discriminative AI** to **generative AI** …

[1]  JHEP05 (2019) 036

# Generative Adversarial Networks (**GAN**)

A generative algorithm, based on an architecture with 2 NNs, a generator G and a discriminator D, which <u>compete</u>

- G creates images from noise, D classifies them real vs fake

- Once trained one against the other, G pursues its goal which is to confuse D, and in the process it learns how to creare fake but very realistic images

**Industry**:
image editing, data generation, security, ..



**HEP**:
Simulate the detector response at reduced computational costs

# Data sparsity and point clouds

HEP handles **high sparsity** datasets: not a HEP-only issue..

- Granularity and occupancy in HEP sensors

- Popolation of stars and galaxies in the cosmos

- Molecular description in computational chemistry

Abstract space with coordinate of sparse elements, each characterized by an array of features, a set of arrays as a function of event/run, …

- e.g. EM shower → E deposited in active volumes of an ECAL

More adequate representation would be a "**point cloud**", and best approach might be not (HEP-)traditional

- Need to be open-minded towards methods not familiar (co far) to the HEP community

- e.g. problems configurable in extraordinarily similar way to how social media datasets are treated (!)



3D semplification

# Graph Neural Network (**GNN**)

Think of a CNN acting on its input features (pixels). Its power resides also on a "regular-array dataset" paradigm

- Data represented as sets of <u>dense</u> arrays/tensors, with intrinsic metrics

In a sparse representation, we need a metric that defines proximity in the abstract space of features

How? Migrating from"**datasets**" to "**graphs**"

- Connect elements of a dataset and train a NN to learn which are the relevant connections

GNN → build a data structure, (V,E) with V=vertex and E=edges, choose possible types of vertices (if no prior one builds a fully-connected graph), etc

GNNs start to be part of present/future of DL in HEP

# ML/DL in HEP

"Traditional" ML   Seeking DL solutions        What next?

# Natural Language Processing (NLP)

Important advances in the decade 2010-2020, thanks to deepNN ("Deep Learning")

- **Machine Translation** → e.g. language A - Language B translation

- **Text/Document Classification** → e.g. doc clustering, sentiment analysis, ..

- **Entity Extraction** (a.k.a. Named Entity Recognition) → extracting relevant information from unstructured text (e.g. vital parameters of a patient from her/his medical record)

- **Summarisation** → e.g. generating concise summaries of docs

- **Question-Answering** → multi-domain and multi-language factual knowledge

- **Digital Assistants** → e.g., Amazon Alexa, Google Assistant, Apple Siri, ..

NLP "first generation" (2010 - ~2017):

- RNN, or Long Short-Term Memory (LSTM), trained on text corpus to learn the structure of language
  - ❖ Handle text as a sequence of symbola ("token"), mapped to multidimensional vectors ("embedding") and processed through a hierarchy of levels

But then.. a "second generation" came in..

*BTW: is NLP relevant for HEP?* → **Yes.**

# "Attention is all you need"



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Feed Forward

Add & Norm

Masked Multi-Head Attention

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Google (Brain), 2017: a new deepNN architecture called **Transformer**

- A native seq-2-seq, with a key element: the "**attention**" mechanism, that allows to pass the meaning of a token in the context of other tokens in the same text

- Soon became the reference model for language processing (see next)



Residual streams: info channel at high capacity (d=12288 in GPT-3). Attention and MLP modules can read/write from/to subspaces of residual streams w/o interference with message passing.

The residual stream is high dimensional, and can be divided into different subspaces.

Layers can interact by writing to and reading from the same or overlapping subspaces. If they write to and read from disjoint subspaces, they won't interact. Typically the spaces only partially overlap.

Layers can delete information from the residual stream by reading in a subspace and then writing the negative verison.

# Language Models: from LM to **L**LM

"Second generation" → modern **Language Models** (**LM**): a class of probabilistic models that learn patterns in NL via more advanced methods

- Architetture change: from RNN to **Transformer**

- Tokenization → multi-language and multi-domain

- Training on extremely large text corpuses (incl. source code!)

Major players:

- **GPT-3** (OpenAI, 2020) from which **ChatGPT** derived (OpenAI, 2022): trained on 45 TB of text (equiv. 2000x Wikipedia), Estimated training cost: 4.6 M$. **GPT-4** (OpenAI, 2023). Estimate: 1 order of magnitude more params

- **Bard/Gemini** (Google), **Claude** (Anthropic), **LLaMA** (Meta)

Significant efforts into scaling LMs into **Large LMs** (**LLM**)

Before getting to this, some examples..

- → training bigger models on more data with greater compute

- → steady+predictable improvements in their ability to learn patterns

- This could be observed in improvements to *quantitative* metrics.. but also *qualitative* (!)

Observations in astronomy:

# LLMs for **Multimessenger Astronomy**

*Description: Referred to by ATel #: 8706, 8718, 8783, 8789 On Jan 14, 2016, the Large Area Telescope (LAT) observed strong* gamma-ray emission *from a new source. The best-fit location of this* gamma-ray source *(RA=8.91 deg, Dec=61.52 deg, J2000.0) has a 95% containment radius of 0.08 deg. This source is not in any published LAT atalog and in the past has not been detected by AGILE or EGRET. The closest candidate counterpart is the* radio source 87GB 003232.7+611352 .

An example text based on **ATel** messages ([astronomerstelegram.org](astronomerstelegram.org))

■ object name; ■ the type of the object or physical phenomena; ■ event type

Goal is to build an information extraction system, i.e. recognize a list of predefined concepts (celestial objects, astronomical facilities, physical properties, people, organization etc.) from a text and produce LLM-generated event summaries based on the parameters of each event

A Cycle 6 ALMA proposal for 4 epochs of stand-alone Atacama Compact Array Band 7 observations of variables in Serpens identified by the JCMT Transient Survey has been accepted (PI: Logan Francis, project code 2018.1.00917.S ). These observations will complement results from the contemporaneous Transient Survey by observing at 850 \u03bc m with a resolution of 38 (compared to the 146 resolution of the JCMT), sufficient to reach the scale of the inner envelopes (\u223c1500 au) of protostars in Serpens.

From entity recognition…

… to semantic relationship extraction

MASTER OT J105440.86-391319.0 - PSN in PGC600519 ( 2.9"E,7"N ) MASTER-SAAO auto-detection system ( Lipunov et al. , " MASTER Global Robotic Net " , Advances in Astronomy , 2010 , 30L ) discovered OT source at ( RA , Dec ) = 10h 54 m 40.86s , -39d 13 m 19s.0 on 2016 - 10 - 08 02:56:46.019UT with m_OT=16.5 . This PSN is in 2.9"E,7"N from the center of PGC600519 and is seen on 8 images . We have reference image without OT on 2015 - 02 - 27 22:35:30UT with m_lim=19.8 Spectral observations are required The discovery and reference images are http://master.sai.msu.ru/static/OT/MASTEROTJ105440.86-391319.0.jpg. MASTER OT J160918.30 - 333456.7 discovery - bright possible dwarf nova outburst , ampl&gt;5.4 m MASTER-SAAO

D. Bonacorsi

*"Cosmosage"*, a general-purpose AI-assistant specialised in answering questions about cosmology (based on Mistral-7B-v0.1)

- training dataset: arXiv papers, astro textbooks, physics textbooks, wikipedia

**user:**
What came before the Big Bang?                                    Popular science

**cosmosage:**
The concept of what came before the Big Bang is a topic of debate and speculation in the realm of cosmology. The Big Bang theory states that our universe began approximately 13.8 billion years ago from a singularity, an infinitely dense and hot point. At this point, space, time, matter, and energy were all created simultaneously. Scientists cannot explain with any certainty why the Big Bang happened, so any speculation about what came "before" is pure theoretical musings. The Big Bang is widely accepted as accounting for everything we have been able to observe about the universe, but it doesn't explain itself. There are many theories and models that attempt to explain the pre-Big Bang universe, such as the cyclic model, the ekpyrotic model, and the string theory landscape. However, these theories are still under development and have not been proven to be accurate descriptions of the actual events that took place before the Big Bang.
**ChatGPT 3.5 gives a similar quality answer.**

**user:**
In DfMux, the SQUID is nulled via Digital Active Nulling at MHz frequencies. How is the SQUID kept nulled at low frequency?                                    Detailed Instrumentation

**cosmosage:**
In DfMux, the SQUID is kept nulled at low frequency by using a low-frequency flux-locked loop (FLL). The FLL is a feedback system that maintains the SQUID in its nulled state by adjusting the input current to cancel out the flux signal. This is necessary because the SQUID is sensitive to magnetic fields, and the magnetic field generated by the input current can change the bias flux, causing the SQUID to lose its null state. The FLL helps to maintain the null state by continuously adjusting the input current to counteract any changes in the bias flux.
**ChatGPT 3.5 says [... passive nulling ... without requiring active intervention...] which is exactly the wrong answer.**

**user:**
Why does a cross-spectrum not suffer from noise bias?                    Data analysis principles

**cosmosage:**
An auto spectrum suffers from noise bias because the noise is correlated with itself, leading to an overestimation of the signal. In contrast, a cross spectrum does not have this noise bias because the noise in different maps is uncorrelated. Therefore, cross spectra are often preferred for cosmological analysis because they provide an unbiased estimate of the signal.
**ChatGPT 3.5 says [... cross-spectra are less prone to noise due to cancellation effects ...] which is incorrect.**

**user:**
What is the current best measure of the optical depth to reionization?                                    Current cosmological parameter constraints

**cosmosage:**
The current best measure of the optical depth to reionization is $\tau = 0.054 \pm 0.007$, as determined by the Planck Collaboration et al. (2018). This value provides strong evidence for reionization at high redshift and helps constrain models of cosmic reionization.
**ChatGPT 3.5 talks about future measurements but does not mention a numerical value or even mentions Planck.**

# A LLM-based AI-assistant for a CERN experiment

"ChATLAS" a **prototype LLM project in a LHC experiment (ATLAS) at CERN** (as of end 2023)

Data **gathering** part is interesting (data **chunking** and data **retrieval** not described here)

- Docs: twiki (>2k), sw docs (>500), e-groups/mails archive (>10k), indico meetings' agendas incl. attached slides and minutes (>440k), Mattermost, Jira tickets, experiment' papers and internal notes (>66k)

- Either HTML or scraped into markdown

Many open challenges:

- highly heterogeneous data

- ensure that collaboration DBs are accessible and exportable; websites should live on a git repo; pubs should be saved as latex, and compiled separately; discussion forums should have anonymisation options… Estimates indicate that this would have saved ~1 yr of data wrangling

- Hallucinations are a real problem

- Not many gpu-hrs, but many expert-hrs, needed for any high-quality fine-tuned AI assistant

# Educational Outreach with **AI-Assisted CERN Open Data**



https://opendata.cern.ch/

D. Bonacorsi

# An example: a Higgs analysis guide

**Define a training goal**

**ATLAS Open Data Higgs analysis guide**

By Philipp Gadow

Guides on Higgs boson data analysis with ATLAS Open Data.

What is the Higgs boson and how can I measure its m...

Show me how to plot the invariant mass distribution ...

How do I read ROOT files for my analysis?

How can I fit the mass peak of the Higgs and how can...

Message ATLAS Open Data Higgs analysis guide...

ChatGPT can make mistakes. Consider checking important information.

**Get draft code for it**

```
sorted_indices = np.argsort(event_E)[::-1]
if len(sorted_indices) < 2:  # Check if there are at least two photons
    continue  # Skip events with fewer than two photons
idx1, idx2 = sorted_indices[:2]  # Indices of the two photons with largest energy

# Extract the leading and subleading photon four-momentum components
E1, E2 = event_E[idx1], event_E[idx2]
eta1, eta2 = event_eta[idx1], event_eta[idx2]
phi1, phi2 = event_phi[idx1], event_phi[idx2]

# Convert eta, phi, E to px, py, pz
px1,
```

Message ATLAS Open Data Higgs analysis guide...

ChatGPT can make mistakes. Consider checking important information.

**Get explanations, refine, learn**

you are calculating the invariant mass using these ene[...]
GeV (gigaelectronvolts), this could be the source of discrepancy leading to an unexpected range or empty plots.

To fix this, you should convert the energies from MeV to GeV before performing the invariant mass calculation. This can be done by dividing the energy values by 1000:

```python
E1, E2 = event_E[idx1] / 1000., event_E[idx2] / 1000.  # Convert from MeV to GeV
```

Incorporate this conversion into your code where you define `E1` and `E2` and then plot the invariant mass distribution again. This adjustment should help ●

Message ATLAS Open Data Higgs analysis guide...

ChatGPT can make mistakes. Consider checking important information.

**Reach the training objective**

Invariant Mass Distribution of the Two Highest Energy Photons

# LLM for **particle accelerators**

A variety of projects..

Plenty of work in progress on LLMs, showing potential towards natural language driven autonomous particle accelerators

- Attempts with GPT 3.5 Turbo, Megadolphin, Vicuna 7B 16K, Mistral 7B, Mixtral 7x8B, Starling-LM, GPT 4 Turbo, GPT4, Orca 2 7B, Orca 2 13B , Llama 2 70B, Falcon 180B, ..

- Constant seek for (and tests with) better models, better prompting, …

GAIA (@**DESY**): a General AI-assistant for Intelligent Accelerator Ops

- Experimental "procedures" defined as a collection of high-level "actions" in a Control System e.g. for managing machine pre-sets

- Exploring a LLM (mixtral:8x7b-instruct-v0.1-q8_0 with 32k context size), agent implemented in Python using the langchain module, prompting based on ReAct (as a combination of chain-of-though prompting and information injection via "actions")
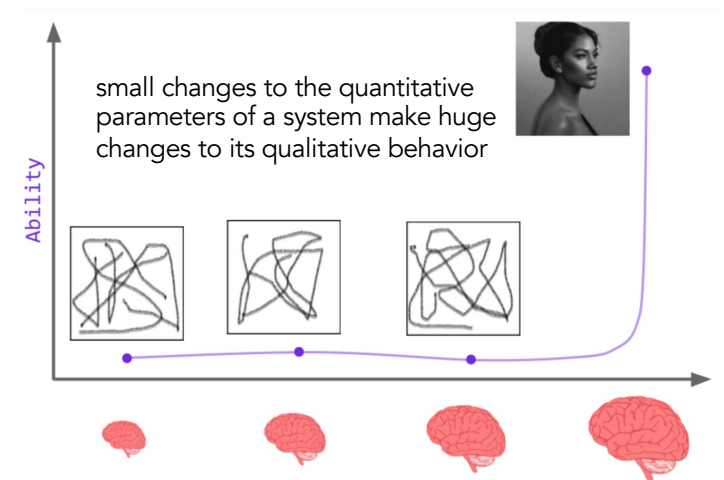
EPA project (@**CERN**), AccGPT, etc..

- EPA = Efficient Particle Accelerator project

- AccGPT = accelerating science via a chatbot for knowledge retrieval for CERN specific content

*[ credits: F. Mayet, J. Kaiser. F. Rehm et al ]*

D. Bonacorsi

# Large Language Models and "emergent abilities"

"**Emergence**": a sudden appearance of a novel behaviour (often referred to as a "phase transition")

- Scaling LLMs → hit a series of critical scales at which new abilities are suddenly "unlocked"

- not directly trained to gain such abilities: they just manifest rapidly and in unpredictable ways
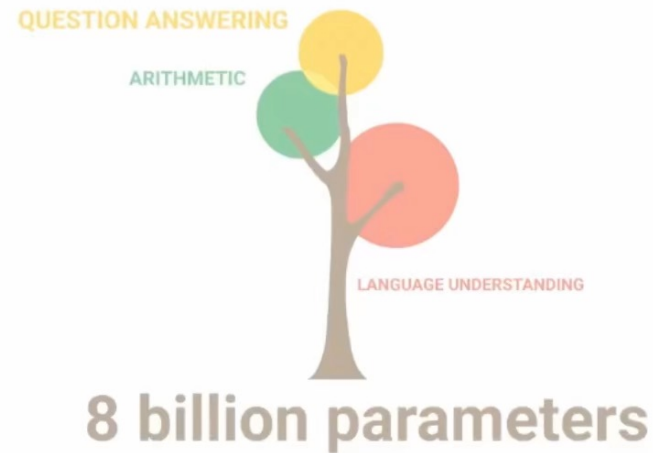
Examples:

- problem solving (math, logic, quantitative reasoning), common sense and social behaviour, (controlled) generation of texts, images, sounds, ..; ability to write, correct, and execute (pseudo)-code



small changes to the quantitative parameters of a system make huge changes to its qualitative behavior

D. Bonacorsi

QUESTION ANSWERING

ARITHMETIC

LANGUAGE UNDERSTANDING

8 billion parameters

# Emergence and prompting

**LLMs can be prompted**: no fine-tuning of a model, but give it NL instructions

**in-context learning**: after training (frozen weights), a model is prompted a set of "gold-standard" examples to illustrate how to complete a task for which it was not trained, and it learns

## Few-shot prompting

emergence at ~100B params for a wide range of cases

→ *prompt engineering, Langchain, ..*

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1   Translate English to French:        ← task description
2   cheese =>                            ← prompt
```

**One-shot**

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1   Translate English to French:        ← task description
2   sea otter => loutre de mer          ← example
3   cheese =>                           ← prompt
```

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1   Translate English to French:        ← task description
2   sea otter => loutre de mer          
3   peppermint => menthe poivrée        ← examples
4   plush girafe => girafe peluche      
5   cheese =>                           ← prompt
```

In addition to in-context learning, another set of interesting abilities stem from **prompt augmentation**

## Chain-of-thought (CoT) reasoning

**Standard Prompting**

*Model Input*

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

*Model Output*

A: The answer is 27. ✗

**Chain-of-Thought Prompting**

*Model Input*

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

*Model Output*

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✓

E.g. it has been demonstrated (in some cases) that just adding "*let's think step by step*" may trigger multi-step reasoning and lead to impressively increased accuracy in arithmetic tasks

D. Bonacorsi

# Biological vs Machine "intelligence"

## Mammalian biological brains

Cat brain

0.760 billion neurons
10 trillion synapses

Human brain

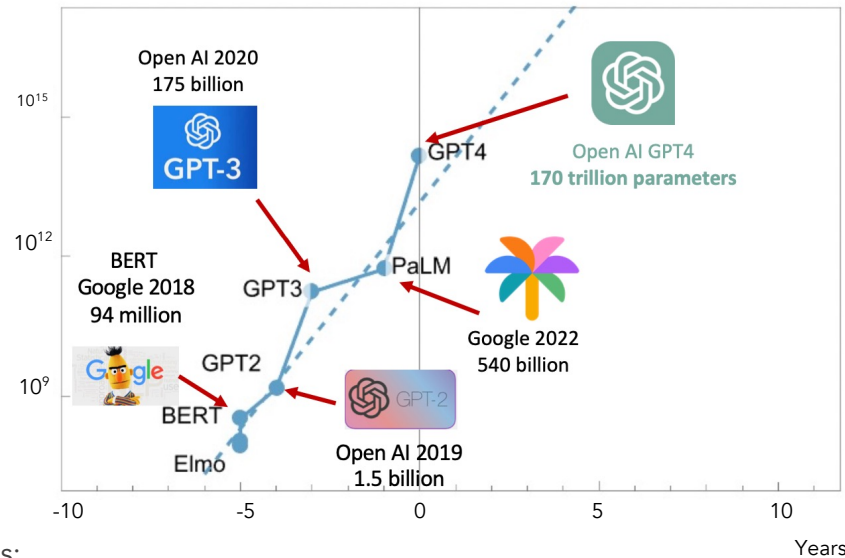80 billion neurons
150 trillion synapses

size of GPT 3.5

size of GPT 4

## LLMs

Parameters



Open AI 2020
175 billion
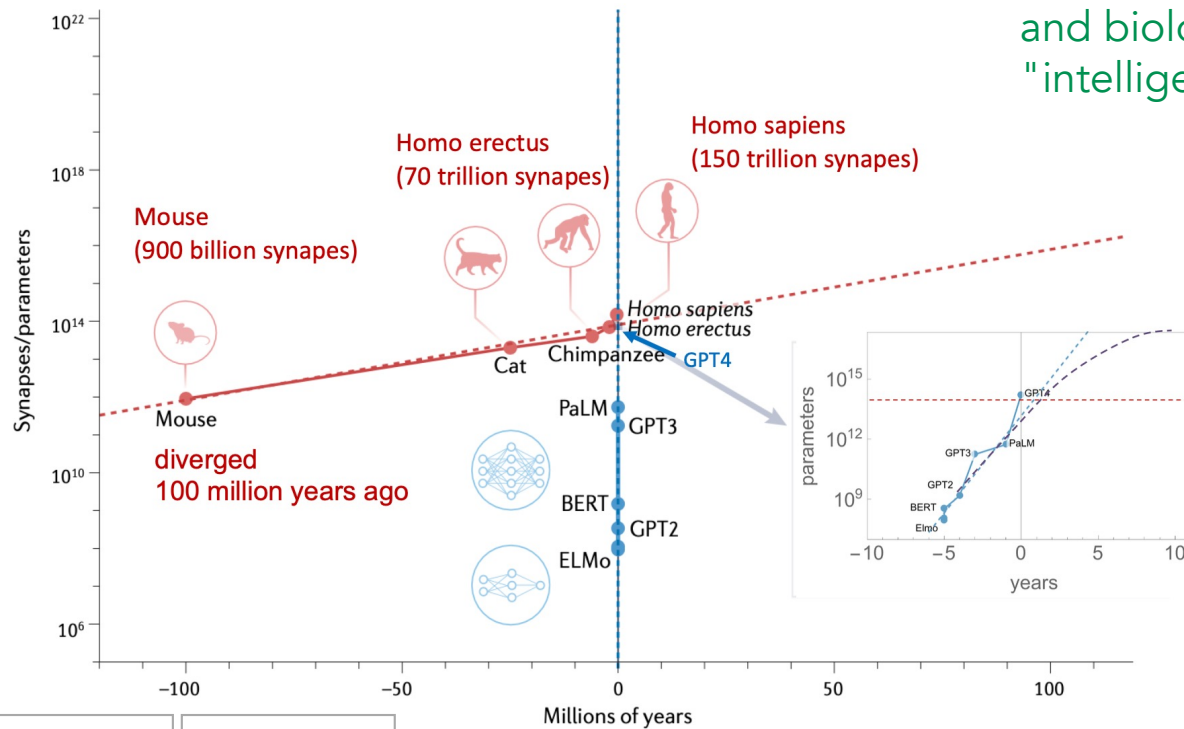
GPT-3

BERT
Google 2018
94 million

GPT3

GPT4

Open AI GPT4
**170 trillion parameters**

PaLM

Google 2022
540 billion

GPT2

BERT

Open AI 2019
1.5 billion

Elmo

Years

Current LLMs:

- Parameters: roughly **the same nb** ($10^{14}$) as the human brain
- .. but **more compute**: brain ($10^{16}$ FLOPS) over a lifetime (100 years) → $10^{22}$ ops, to be compared with LLM training time, around $10^{25}$ ops
  - ❖ And it consumes more.. **Red AI** is a serious issue!

# Biological vs Machine "intelligence"

Measuring "intelligence" by number of neurons (or computational units):

- biological → growth by a factor 2x in 1 million years

- machine → growth by a factor 10x in 1 year

The intersection - when machines and biology have comparable "intelligence" - is **~now**
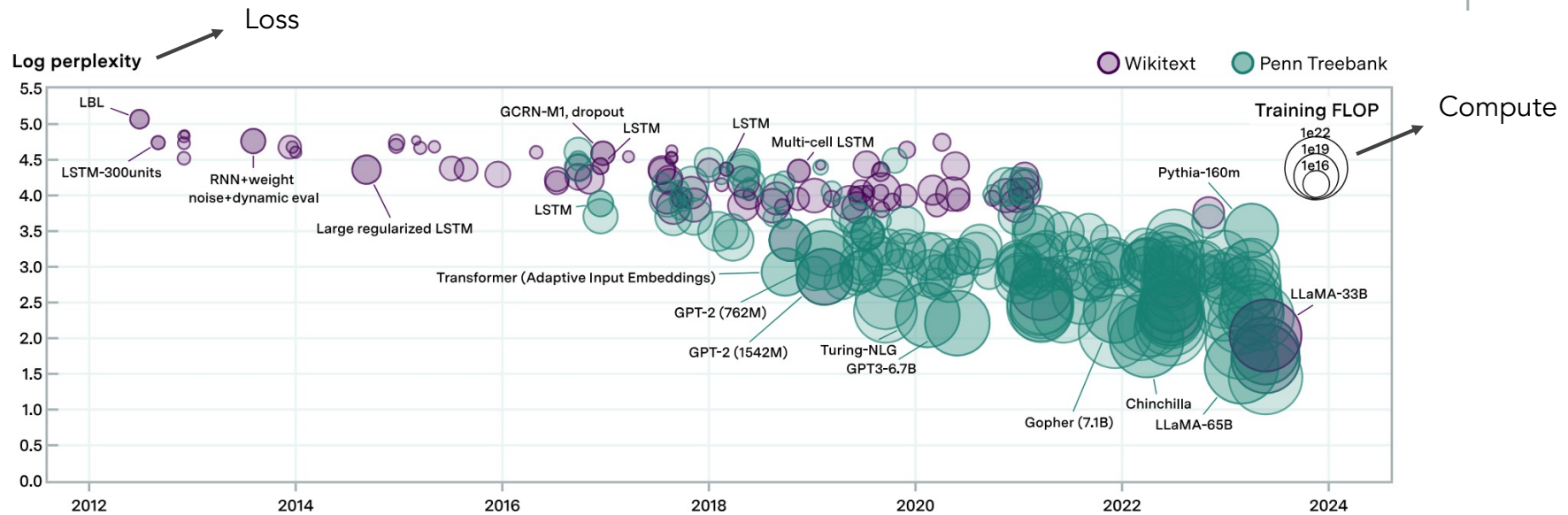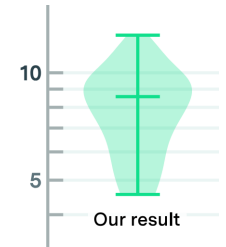


even sub-exponential grown will soon be superhuman !

| Credits: M. Schwartz

47

D. Bonacorsi

# Algorithmic progress in LLMs

The compute required to reach a set performance threshold has **halved approx. every 8 months**, with a 95% confidence interval of around 5 to 14 months

- Algorithmic improvements **faster** than hardware gains per Moore's Law!





Will we run out of data? compute? networks? … or **energy**?

.. and algorithms will continue to get better → also if/when **written by AI** itself..
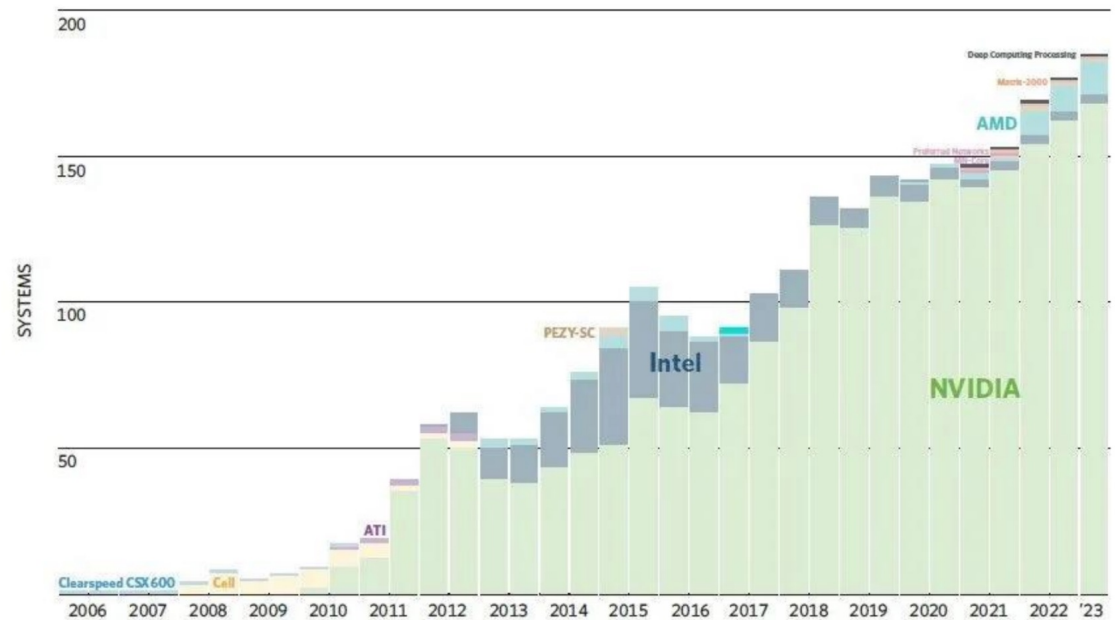
D. Bonacorsi

# A glance to hardware: the NVIDIA "gravity"

Of the 184 accelerated machines on the TOP500 (June 2023) list, 167 have Nvidia GPUs

- 11 of them have AMD GPUs

- 6 have other kinds of accelerators like Intel Knights coprocessors.

AMD has 5.9 percent of accelerated system share in GPUs, compared to 90.8 percent for Nvidia (95% in AI)
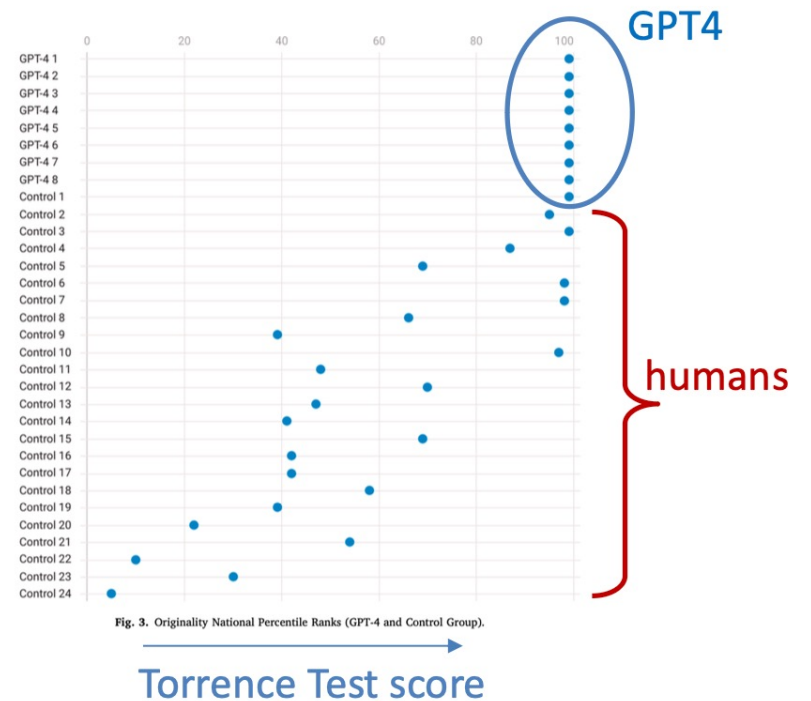
Nvidia is <u>not</u> slowing down in 2024..

# "but physics requires creativity.." (?!)

Article | Open access | Published: 10 February 2024

## The current state of artificial intelligence generative language models is more creative than humans on divergent thinking tasks

Kent F. Hubert ✉, Kim N. Awa & Darya L. Zabelina

The emergence of publicly accessible artificial intelligence (AI) large language models such as ChatGPT has given rise to global conversations on the implications of AI capabilities. Emergent research on AI has challenged the assumption that creative potential is a uniquely human trait thus, there seems to be a disconnect between human perception versus what AI is objectively capable of creating. Here, we aimed to assess the creative potential of humans in comparison to AI. In the present study, human participants (N = 151) and GPT-4 provided responses for the Alternative Uses Task, Consequences Task, and Divergent Associations Task. We found that AI was robustly more creative along each divergent thinking measurement in comparison to the human counterparts. Specifically, when controlling for fluency of responses, AI was more original and elaborate. The present findings suggest that the current state of AI language models demonstrate higher creative potential than human respondents.
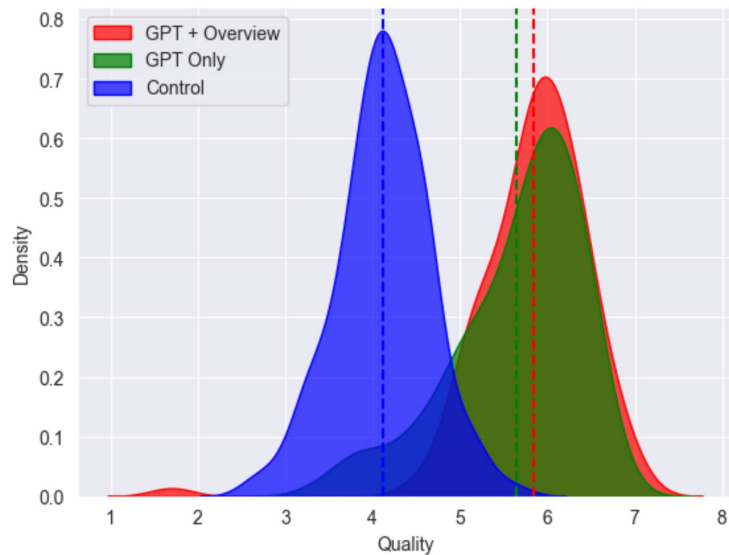
Fig. 3. Originality National Percentile Ranks (GPT-4 and Control Group).

Torrence Test score
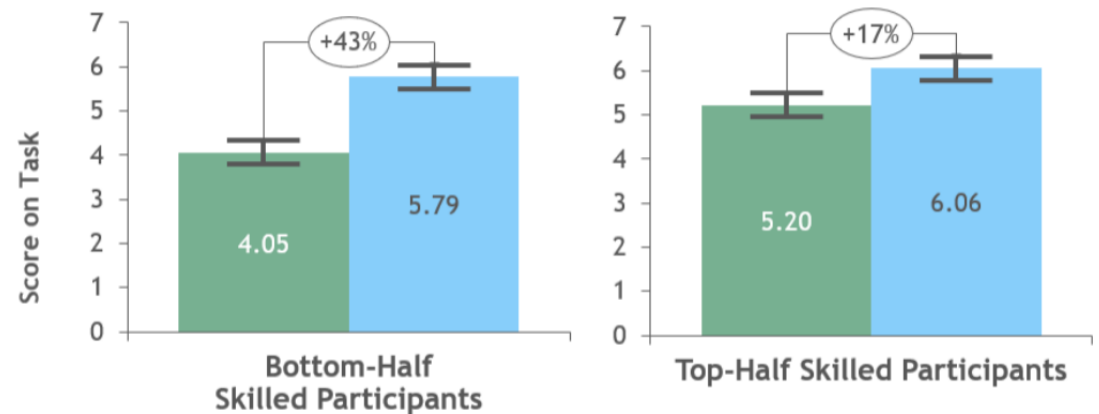
GPT4 more creative than 99% of humans..

# Is this "augmented" intelligence?

## Can LLM help "consultants"?



**Yes, and by a lot!**

## AI-assistants as skill-levellers



Average consultants using AI perform as good as best consultants

Navigating the Jagged Technological Frontier (…), Harvard Business School, 2023)
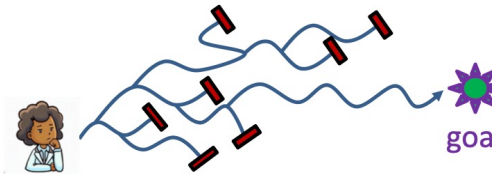
D. Bonacorsi

# Machines not limited by biology

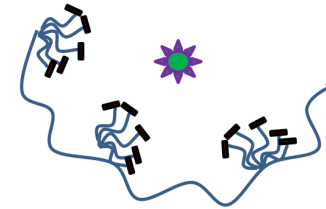## Can AI be a skill-leveller (or more) for e.g. HEP theory?



Credits: M. Schwartz, Harvard and NSF IAIFI, elaborating on hard theoretical physics problems and AI, at EuCAIFCon (Amsterdam, 2024)
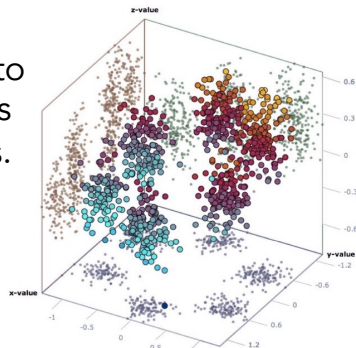
In the past, we made progress depsite many dead ends

Are we even making forward progress anymore?

goal

## Maybe the problems are just **too difficult** (for us)?

- E.g. could a cat ever learn to play chess? Humans have limits too.. i.e. biology

Humans like to "visualise", as we have eyes.



→ project in 2D

For a machine, 2D is not special: it can easily visualise in d dimensions

Humans hold few concepts in working memory at once, and like "simple and elegant" equations

$$i\partial_t \psi = H\psi$$
$$G_{\mu\nu} = \kappa T_{\mu\nu}$$

A computer memory can handle much more concepts at once, and can understand systems not governed by simple equations

D. Bonacorsi

# We are a training set for machines

Current state-of-the-art AI can answer questions / (~) solve textbook problems

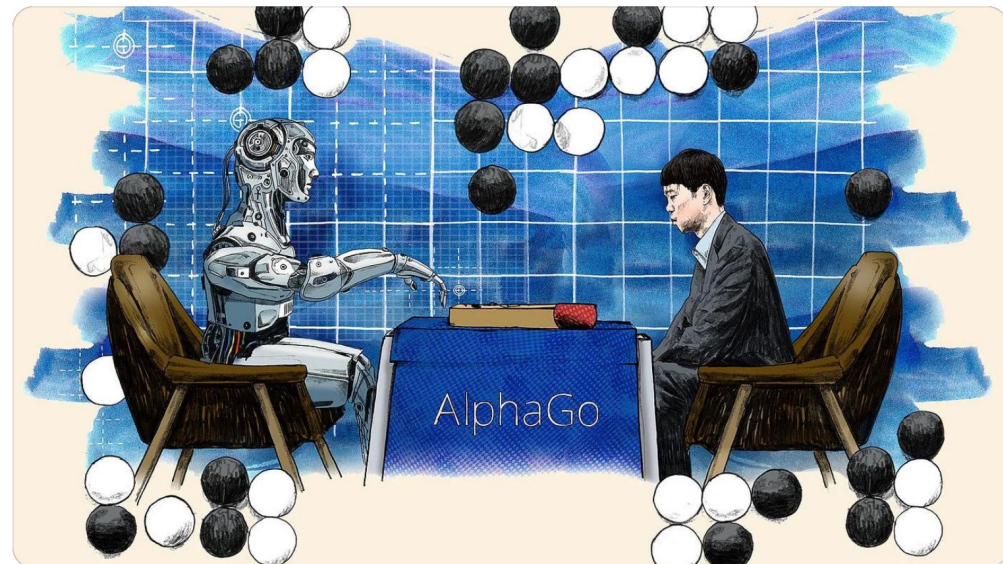How? → via training on huge datasets of answered questions / solved problems

By whom? → **Us**! **we** answered and solved all that, **we actually generated its training set**

- (and we do the same for ourselves)

E.g. LLMs:

- learn from our training set

- Human feedback helps refine the models

- Machines generates data and refines its models

Humans and machines seem very close to be not so different..
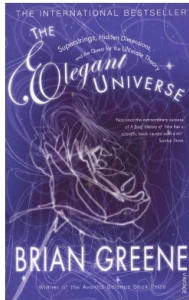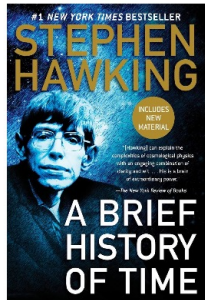


53

D. Bonacorsi

# (Beyond) augmented intelligence

Suppose a machine understands the theory of everything but we don't
- e.g. can calculate the fine-structure constant from scratch
- e.g. can preduct the endpoint of black-hole evaporation

**Is this enough or do we need to understand it too?**

- **The authors of Popular science books** understand the details; we just get the general idea

I don't understand the proof of Fermat's last theorem
- I'm glad that somebody does
- Does it matter that the person is human?

If a machine understands fundamental physics it can
1. Dumb it down so we can get the general idea
2. Find practical applications

Is this what we want? **I guess not.**

But..

What if this is the best we will get? What if AI makes us optimistic for substantive progress in HEP theory in our lifetime?

D. Bonacorsi

# So, in a nutshell: <u>is AI all we need?</u>

AI is not "all" we need, but "something" that we need for sure, at some level.

- Software/Computing challenges to keep up with HEP goals are tough: AI as "**part of the solution**"

Nevertheless, we had better start envisioning ML/DL not as mere "tools"

- Think of this as a **discipline** that showed up to be impressively useful, but whose theory is "under construction"

Our **best**:

- **Early career researchers**

- **Domain knowledge** and scientific rigour, ability to explore what's solid and ignore the hype, direct efforts.

Our **risk**:

- **Hard to keep the pace**, but reluctant (time-wise, funding-wise, ..) to admit that this is <u>not</u> a part-time job

- R&D vs production: need more efforts to bring approaches to **production quality** (e.g. DL pipelines on AI infrastructures)

Given the AI recent past, in a 5-years timescale something radically new may probably come up.

**If we fail to prepare** for it, let's **prepare to fail** in exploiting its value