# Towards AI-driven User Support for CNAF Tier-1

**M. Barbetti**, A. Trashaj, E. Ronchieri, C. Pellegrino, and the contribution of many others

**CNAF seminar** | 21 November 2024

**1.0**

# THE EVOLUTION OF ARTIFICIAL INTELLIGENCE

*A sensationalistic overview of the history of Artificial Intelligence*

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# *Can you find the (AI)ntruder?*



**Credits:** Jorg Karg



**Credits:** Boris Eldagsen



**Credits:** Aleksei Andreev



**Credits:** Evgeny Lutsko

3

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing        Missione 4 • **Istruzione e Ricerca**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# Can you find the (AI)ntruder?

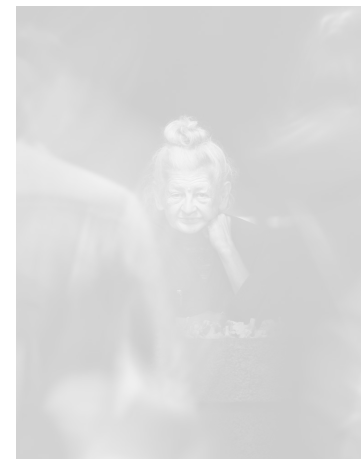**PSEUDOMNESIA: *The Electrician***, 2023



**Credits:** Jorg Karg
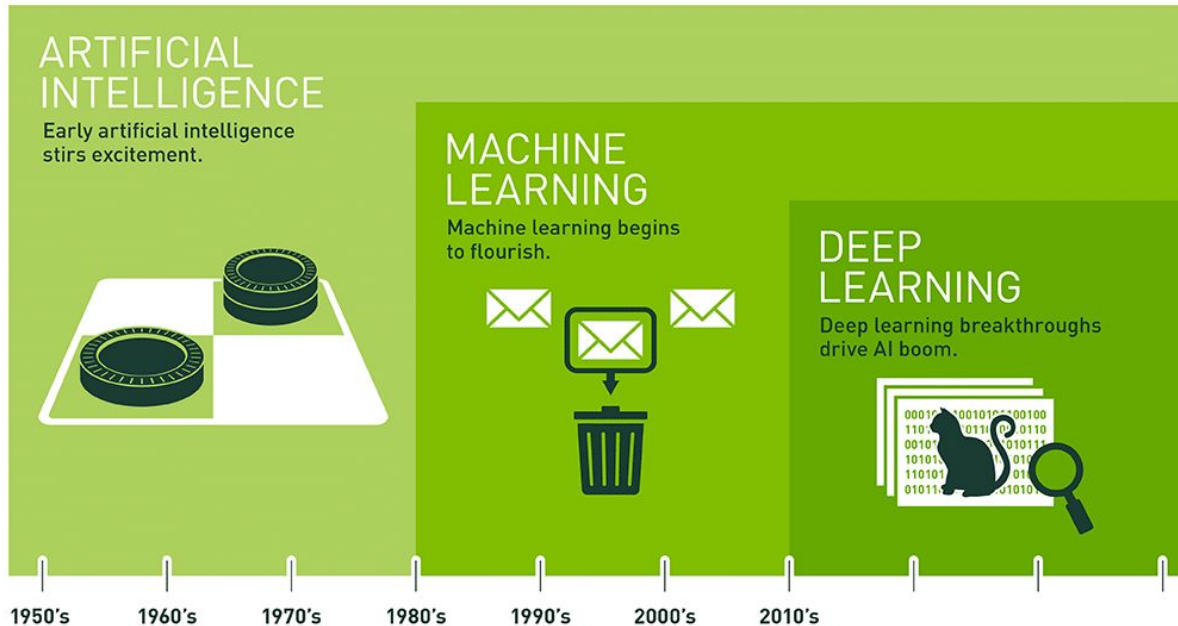


**Credits:** Boris Eldagsen
(by using DALL-E 2)



**WINNER**
of the creative open category at the
*Sony World Photography Awards*
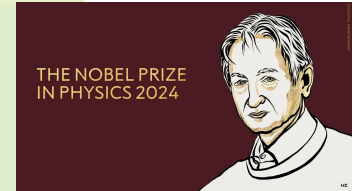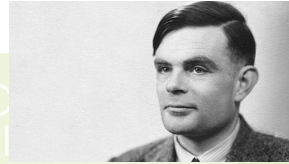
**Credits:** Aleksei Andreev



**Credits:** Evgeny Lutsko

4

# History of Artificial Intelligence

# History of Artificial Intelligence

Alan Turing introduced a test to answer the question "Can machines think?"

First model of an *Artificial Neural Network* by W. McCulloch and W. Pitts

The term *Artificial Intelligence* was coined for a thematic workshop hosted at Dartmouth

THE NOBEL PRIZE IN PHYSICS 2024

Geoffrey Hinton coined the term *Deep Learning* to refer to *deep* ANN-powered algorithms

1960's 1970's 1980's 1990's 2000 10's

**1950**

**1943**

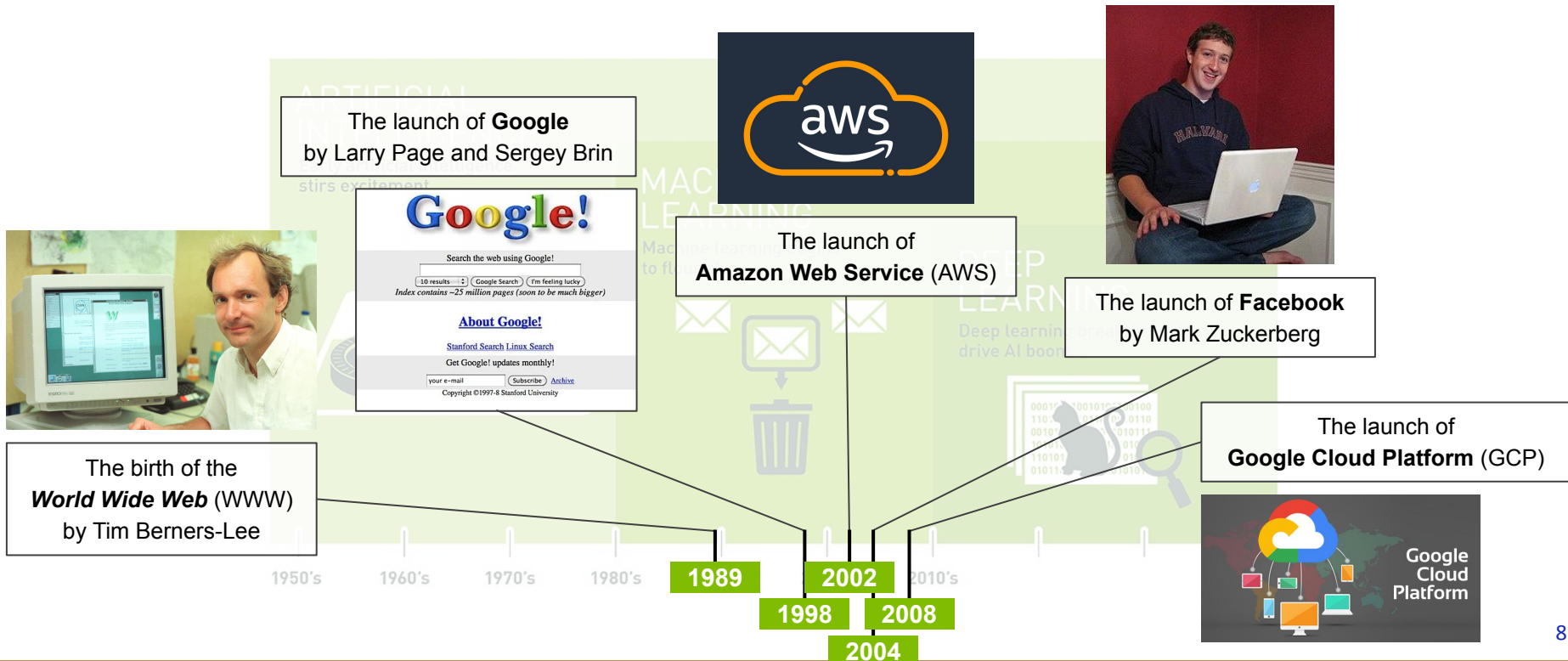**1956**

**2006**

6

# History of Artificial Intelligence



The modern technique for training MLPs, known as *backpropagation*, was first published by S. Linnainmaa
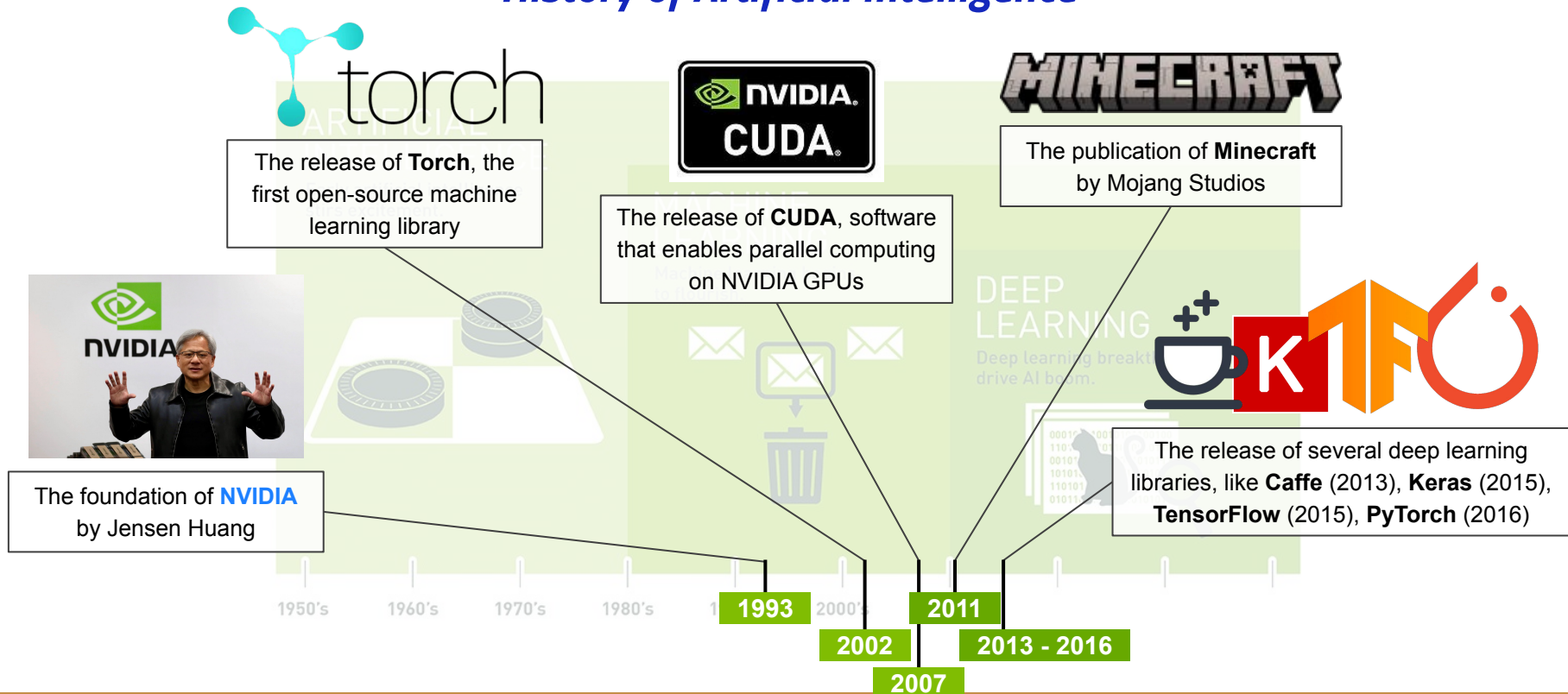
P. Werbos and D.E. Rumelhart independently developed and applied backpropagation to MLPs in the way that has become **standard**

1950's   1960's   **1970**   **1982**   1990's   2000's   2010's
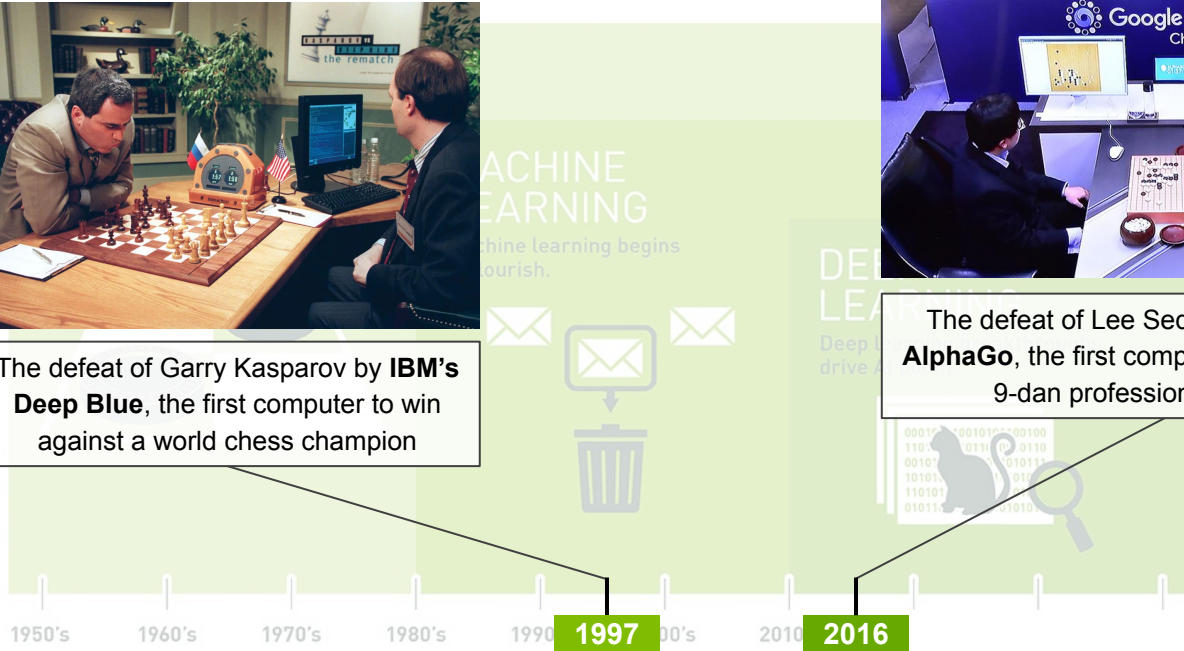
# History of Artificial Intelligence

The launch of **Google** by Larry Page and Sergey Brin

The launch of **Amazon Web Service** (AWS)

The launch of **Facebook** by Mark Zuckerberg

The launch of **Google Cloud Platform** (GCP)

The birth of the **World Wide Web** (WWW) by Tim Berners-Lee

1950's   1960's   1970's   1980's   2010's

**1989**   **2002**

**1998**   **2008**

**2004**

# History of Artificial Intelligence



The release of **Torch**, the first open-source machine learning library

The release of **CUDA**, software that enables parallel computing on NVIDIA GPUs

The publication of **Minecraft** by Mojang Studios

The foundation of **NVIDIA** by Jensen Huang

The release of several deep learning libraries, like **Caffe** (2013), **Keras** (2015), **TensorFlow** (2015), **PyTorch** (2016)

1950's   1960's   1970's   1980's   **1993**   2000's   **2011**

**2002**   **2013 - 2016**

**2007**

# *History of Artificial Intelligence*



The defeat of Garry Kasparov by **IBM's Deep Blue**, the first computer to win against a world chess champion

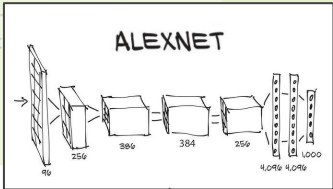The defeat of Lee Sedol by **DeepMind's AlphaGo**, the first computer to win against a 9-dan professional Go player
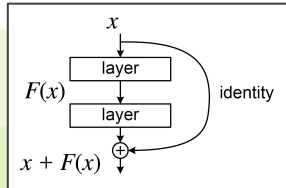
1950's    1960's    1970's    1980's    1990    **1997**    00's    2010    **2016**

10

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing        Missione 4 • **Istruzione e Ricerca**

# History of Artificial Intelligence

The victory of the *ImageNet 2012 Challenge* by a CNN-based model called **AlexNet**

The invention of *Generative Adversarial Nets* by Ian Goodfellow and colleagues

The victory of the *ImageNet 2015 Challenge* by a novel architecture introduced with **ResNet**

The first **ImageNet Large Scale Visual Recognition Challenge** (ILSVRC), an annual contest on image classification and detection via software programs

The launch of **DALL-E** by OpenAI

1950's    1960's    1970's    1980's    1990's    2000's

**2010**    **2015**

**2012**    **2021**

**2014**

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
Big Data and Quantum Computing

# History of Artificial Intelligence



The introduction of **BERT** by Google and **GPT-1** by OpenAI, models for NLP problems

The introduction of the **Transformer** architecture with the paper "Attention is all you need"

The launch of several chatbot platforms, like **ChatGPT** by OpenAI (2022), **Bard/Gemini** by Google (2023), **Claude** by Anthropic (2023)

The release of **ELIZA** by Joseph Weizenbaum, the computer program to engage conversations with humans

ChatGPT

| 1950's | 1960's | 70's | 1980's | 1990's | 2000's | 2010's |
|--------|--------|------|--------|--------|--------|--------|

1966

2017

2018

2022 - 2023

12

# The evolution of image generation

**DCGAN** (arXiv:1511.06434), 2015



**Credits:**

- Alec Radford
- Luke Metz
- Soumith Chintala

**Sora** (OpenAI text-to-video model), 2024



`<PROMPT>` A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about. `</PROMPT>`

13

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing    Missione 4 • **Istruzione e Ricerca**

Finanziato dall'Unione europea
NextGenerationEU

Ministero dell'Università e della Ricerca

Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA

ICSC Centro Nazionale di Ricerca in HPC, Big Data and Quantum Computing

# Chatbot, chatbot everywhere!



**Credits:** https://yellow.ai/blog/future-of-chatbots

**2.0**

# EVOLVING THE CNAF USER SUPPORT THROUGH AI-DRIVEN SOLUTIONS

*How AI technologies can help in supporting Tier-1 users*

Finanziato dall'Unione europea NextGenerationEU

Ministero dell'Università e della Ricerca

Italiadomani
PIANO NAZIONALE DI RIPRESA E RESILIENZA

ICSC
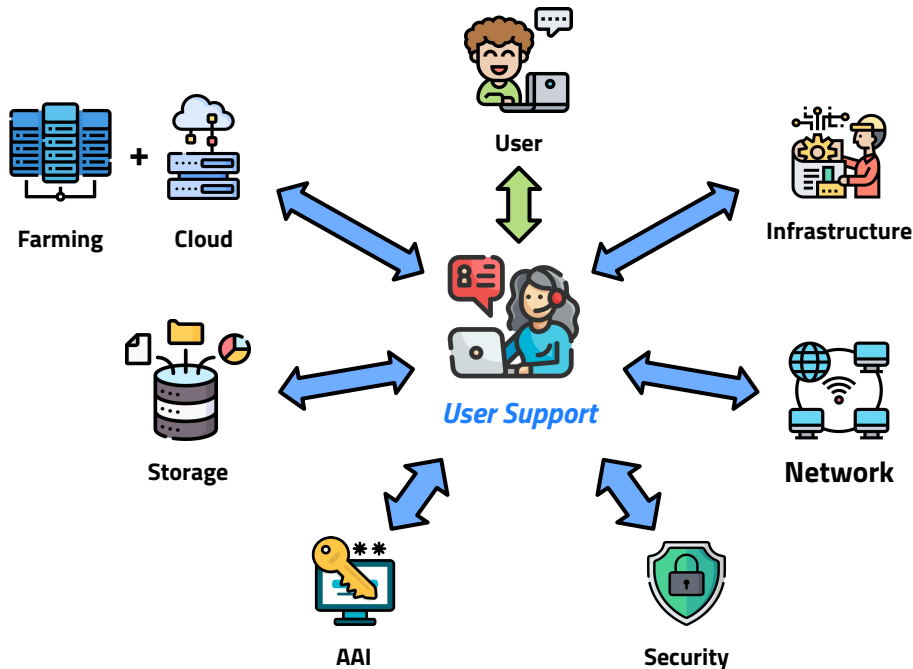Centro Nazionale di Ricerca in HPC, Big Data and Quantum Computing

# *Scientific computing at CNAF*

As INFN national center dedicated to R&D on **Information and Communication Technologies**, CNAF host the Italian Tier-1 data center since 2003, providing computing resources to a **wide scientific community**:

- collaboration with **60+ experiments** (not only the four major experiments of LHC)

- support to **1500+ active users** accessing the Tier-1 resources (including communities different from the High Energy and Nuclear Physics one)

16

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing      Missione 4 • **Istruzione e Ricerca**
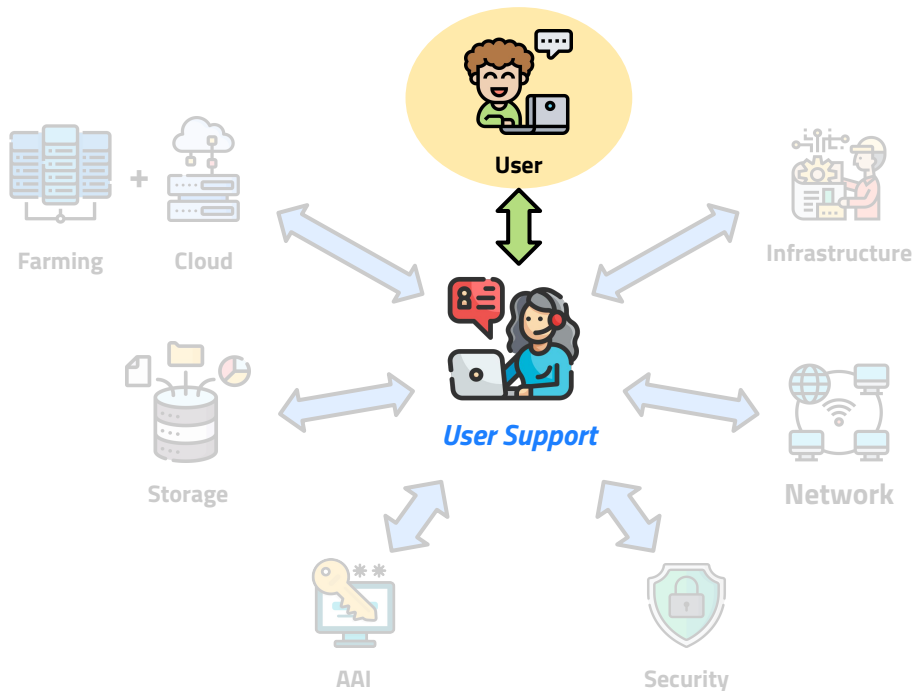
# Support for experiments at Tier-1

The Italian Tier-1 ensures support for experiments and users through the dedicated *User Support* (US) unit:

- **principal contact point** between users and other specialized Tier-1 units
- development of tools/procedures to **simplify** the access to computing resources
- preparation and updates of the **Tier-1 User Guide** (https://l.infn.it/t1guide)

**Three levels of support:**

1. User Support unit (entry point)
2. involvement of other specialized units
3. involvement of software developers



Farming + Cloud

User

Infrastructure

Storage

User Support

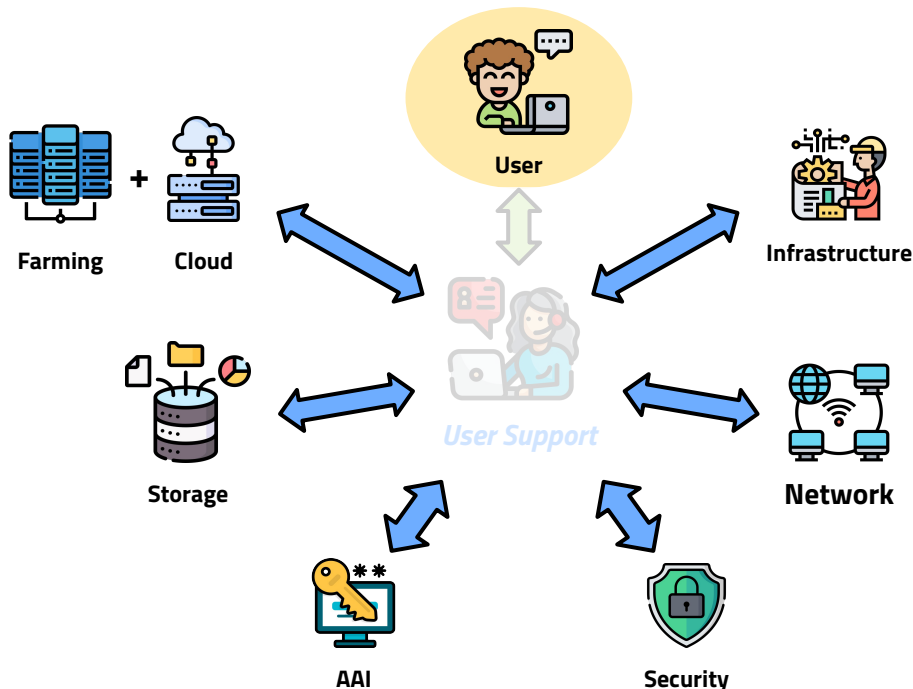Network

AAI

Security

# Support for experiments at Tier-1



The Italian Tier-1 ensures support for experiments and users through the dedicated *User Support* (US) unit:

- **principal contact point** between users and other specialized Tier-1 units
- development of tools/procedures to **simplify** the access to computing resources
- preparation and updates of the **Tier-1 User Guide** (https://l.infn.it/t1guide)

**Three levels of support:**

1. User Support unit (entry point)
2. involvement of other specialized units
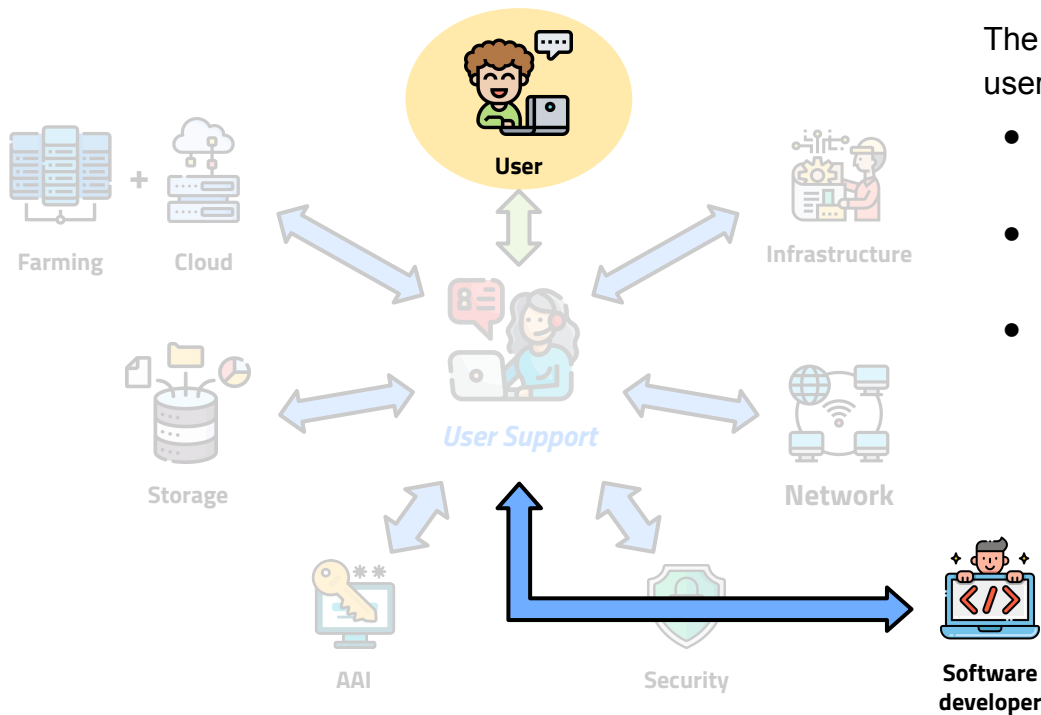3. involvement of software developers

# Support for experiments at Tier-1



The Italian Tier-1 ensures support for experiments and users through the dedicated *User Support* (US) unit:

- **principal contact point** between users and other specialized Tier-1 units

- development of tools/procedures to **simplify** the access to computing resources

- preparation and updates of the **Tier-1 User Guide** (https://l.infn.it/t1guide)

**Three levels of support:**

1. User Support unit (entry point)
2. involvement of other specialized units
3. involvement of software developers
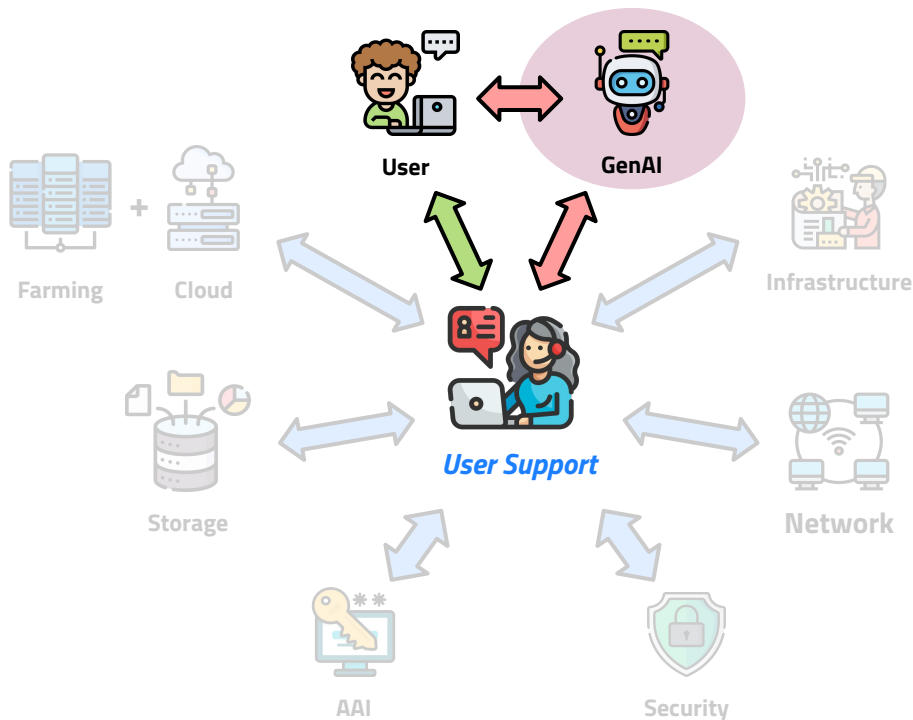
19

# Support for experiments at Tier-1



The Italian Tier-1 ensures support for experiments and users through the dedicated *User Support* (US) unit:

- **principal contact point** between users and other specialized Tier-1 units

- development of tools/procedures to **simplify** the access to computing resources

- preparation and updates of the **Tier-1 User Guide** (https://l.infn.it/t1guide)

**Three levels of support:**

1. User Support unit (entry point)

2. involvement of other specialized units

3. involvement of software developers

20

# *Evolving the User Support using Artificial Intelligence*



**User**

**GenAI**

**Farming**  **Cloud**

**Infrastructure**

**Storage**

**User Support**

**Network**

**AAI**

**Security**

In the context of NRRP and thanks to the action of the ICSC and TeRABIT projects, the **Italian Tier-1 will evolve** significantly in the next few years:

- acquisition of **more and more users and experiments**
- wide scientific community with **various requirements** in terms of infrastructure, hardware, and software
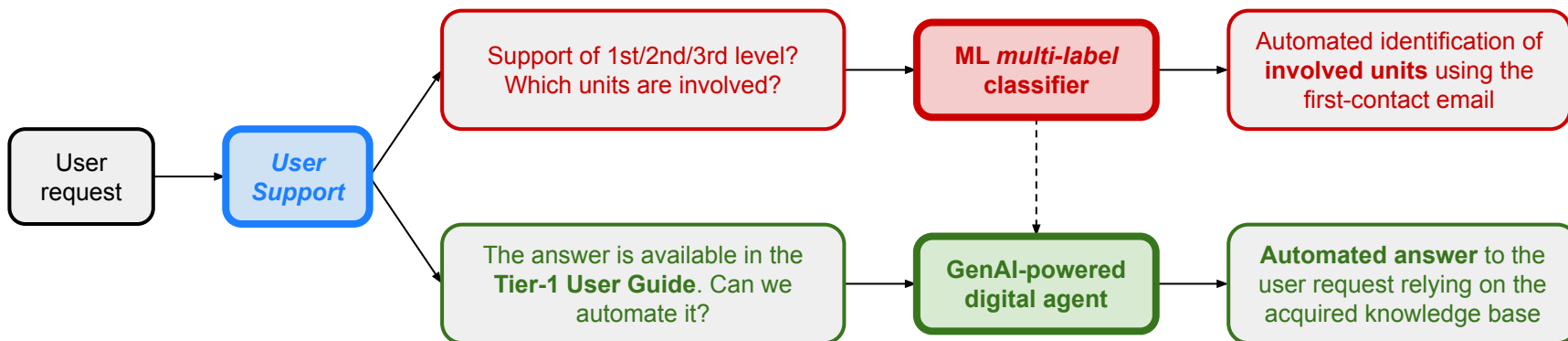- need for an **evolution** of the CNAF data center and its units

**IDEA**

Employing AI models to build a **Digital User Supporter** able to help in supporting a growing number of new users and in adopting the latest software technologies

21

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# Prototyping a Digital User Supporter

The process of developing an **AI-powered digital agent** for managing User Support duties has been organized into two parallel and complementary sub-projects:
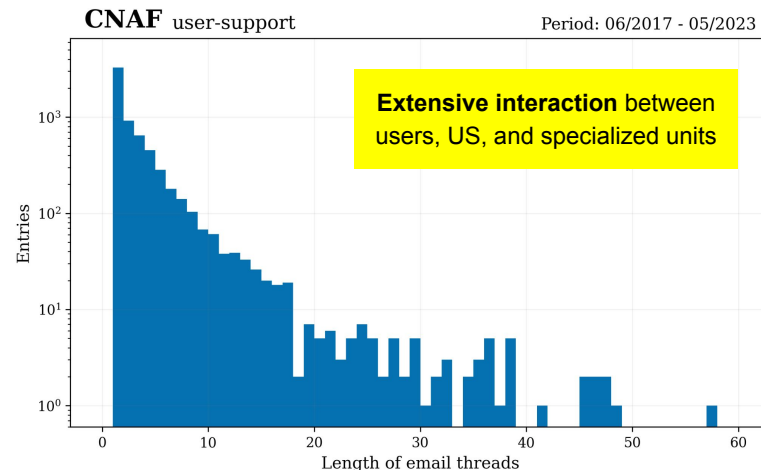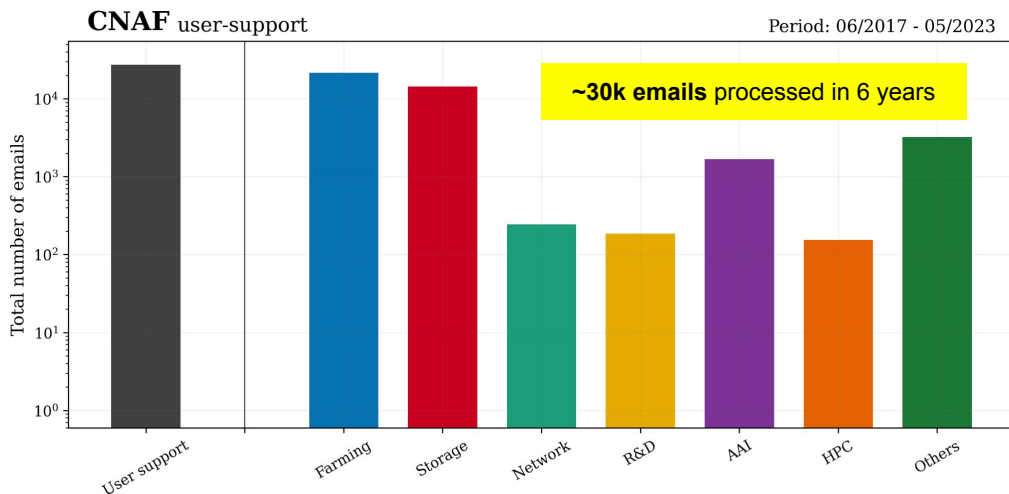
2.1    *Multi-label classification*

# *Steps for building an email classifier*

**1** Preparation of the **dataset of emails** for training

**2** Mapping of the textual features into a **numerical representation** (*embedding*)

**3** **Classifier training** using differente ML algorithms

**4** Benchmark of the **classification performance** and **results combination** (*ensemble*)

24

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing       Missione 4 • **Istruzione e Ricerca**

# A quick look at the data sample



**CNAF** user-support — Period: 06/2017 - 05/2023

~30k emails processed in 6 years

Categories: User support, Farming, Storage, Network, R&D, AAI, HPC, Others

**CNAF** user-support — Period: 06/2017 - 05/2023

**Extensive interaction** between users, US, and specialized units

# Preparation of the training dataset

## Data collection

The sample contains **~30k emails** received/sent in the period 06/2017 – 05/2023

Each has been saved as an individual JSON file (**~260 MB**)

No attachments, and (almost) no HTML content

Each file includes:
- `from` – email address
- `to` – email address
- `date` – datetime
- `subject` – text
- `content` – text
- `parent` – email address

## Anonymization

All the **email addresses** have been **anonymized** by replacing them with UUIDs

Any **references to real individuals** (e.g., names or usernames) found in the email subject/body have been **anonymized** and replaced with the placeholder `[NAME]`

The anonymization process was carried out using **automated scripts** or by **editing manually** the email metadata

## Labelling

Each email is assigned a **label** based on the text contained in its subject/body

The **email thread** between a user and the various Tier-1 units involved is reconstructed using the information in the `parent` field

All labels associated to the same thread are propagated to the first contact email as a part of the **labelling** process
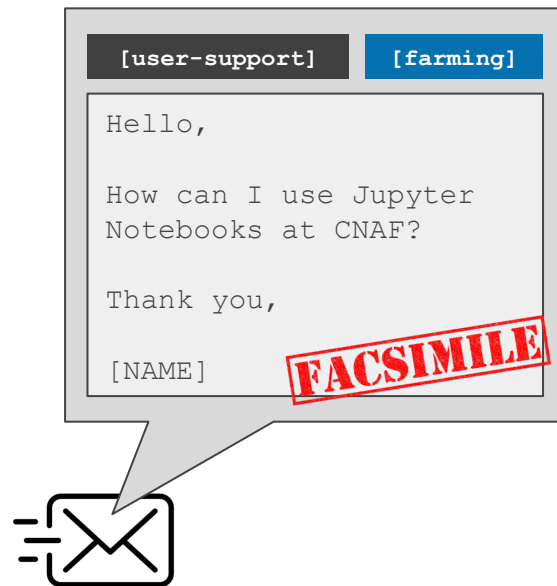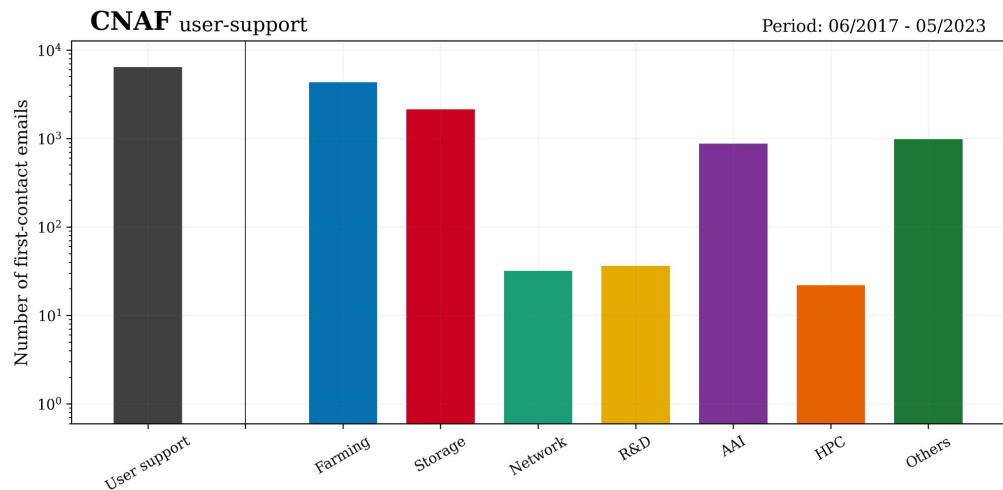
## Data cleaning

The email body may contain **escape sequences** (e.g., `\n`, `\t`) o **HTML tags**

The emails are **multilingual** (English/Italian) and may include **special characters** (e.g., accented letters, emojis)
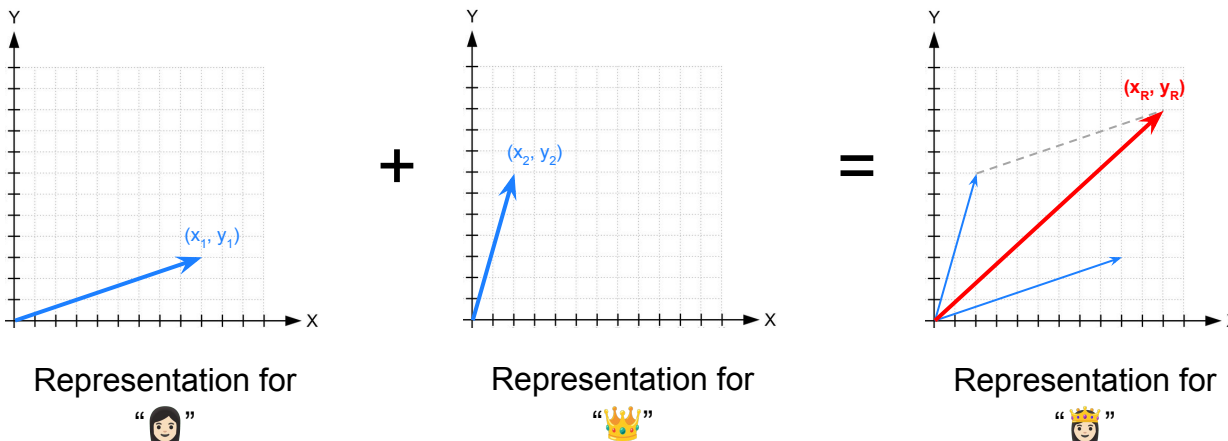
The **data cleaning** process transforms the text to retain only words, placeholder, and punctuation

# The prepared training set

Finanziato
dall'Unione europea
NextGenerationEU

Ministero
dell'Università
e della Ricerca

Italiadomani
PIANO NAZIONALE
DI RIPRESA E RESILIENZA

ICSC
Centro Nazionale di Ricerca in HPC,
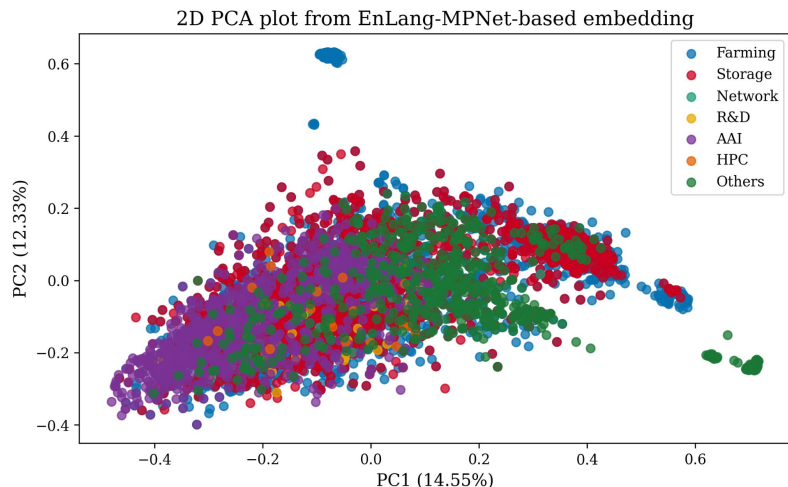Big Data and Quantum Computing

# How to represent the textual features

- The meaning of a text is encoded in word combinations → ML models operate on vectors (**features**)
  - Word (or token) combinations can be mapped into a **"representation space"** (a high-dimensional vector space) through an *embedding model*

- The embedding model can be designed to preserve **semantics**, even across multiple languages
  - **Sentence-Transformers** is a Python framework that provides state-of-the-art models for text embedding



Representation for
"👩🏻"

**+**

Representation for
"👑"

**=**

Representation for
"👸🏻"

28

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# How emails look like in the representation space

2D PCA plot from EnLang-MPNet-based embedding



**all-mpnet-base-v2**

- **base model:** MPNet
- **embedding dim:** 768
- **max seq length:** 384
- **inference speed:** 1
- **language:** English only

**all-MiniLM-L6-v2**

- **base model:** MiniLM
- **embedding dim:** 384
- **max seq length:** 256
- **inference speed:** x5
- **language:** English only

**paraphrase-multilingual-mpnet-base-v2**

- **base model:** XLM-RoBERTa
- **embedding dim:** 768
- **max seq length:** 128
- **inference speed:** ~1
- **language:** multiple

**paraphrase-multilingual-MiniLM-L12-v2**

- **base model:** MiniLM
- **embedding dim:** 384
- **max seq length:** 126
- **inference speed:** x2.5
- **language:** multiple

29

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# *Performance of the trained classifiers*

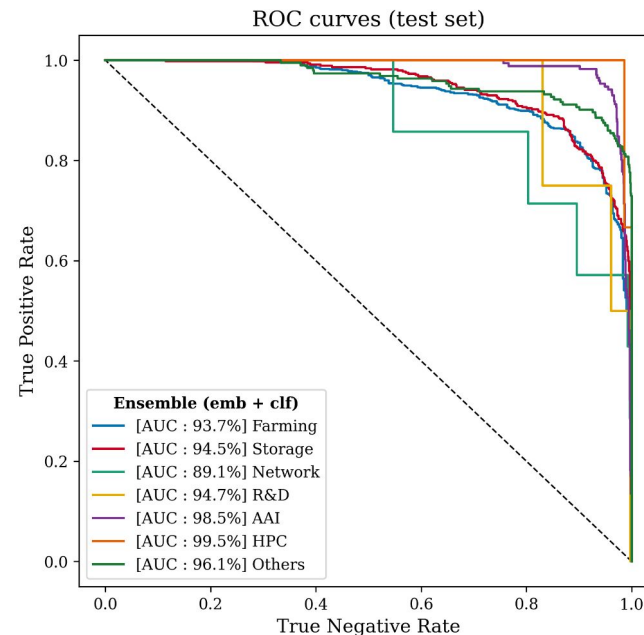Different ML models were trained to perform **multi-label classification**

- ***k-Nearest Neighbors*** (kNN)
- ***Random Forest*** (RF)
- ***Extreme Gradient Boosting*** (XGBoost)
- ***Feed-forward Neural Network*** (FNN)

Various embedding models were used in combination with these classification models

- `all-mpnet-base-v2`
- `all-MiniLM-L6-v2`
- `paraphrase-multilingual-mpnet-base-v2`
- `paraphrase-multilingual-MiniLM-L12-v2`

Different preprocessing strategies → **4 (emb) x 4 (prep) x 4 (clf) = 64 models in total**

By combining the outputs of the most promising models using an **ensemble approach**, the best performance was achieved: **~95% of ROC AUC score** on the test set



ROC curves (test set)

Ensemble (emb + clf)
[AUC : 93.7%] Farming
[AUC : 94.5%] Storage
[AUC : 89.1%] Network
[AUC : 94.7%] R&D
[AUC : 98.5%] AAI
[AUC : 99.5%] HPC
[AUC : 96.1%] Others

True Positive Rate

True Negative Rate

30

2 . 2 *Digital User Supporter*

# The problem of artificial hallucinations



**Credits:** https://fisicisenzapalestra.com/artificiale-come-lintelligenza-e-gli-allucinogeni.html

One of the major challenges in working with AI chatbots and Computer Vision tools is the issue of *artificial hallucinations*

- responses that contain **false or misleading information** presented as **facts**

A potential solution to limit the "creativity" of Large Language Models (LLMs) is to integrate them into an architecture designed to **retrieve relevant information** and compose **answers based on that data**

The *Retrieval-Augmented Generation* (RAG) architecture provides a **semantic-aware pipeline** that retrieves information from a knowledge base and generates text grounded in the retrieved data

**PROTOTYPE**

The first prototype of the **AI-powered Digital User Supporter** has been developed upon a RAG model, implemented with LangChain and provisioning the LLM relying on ollama
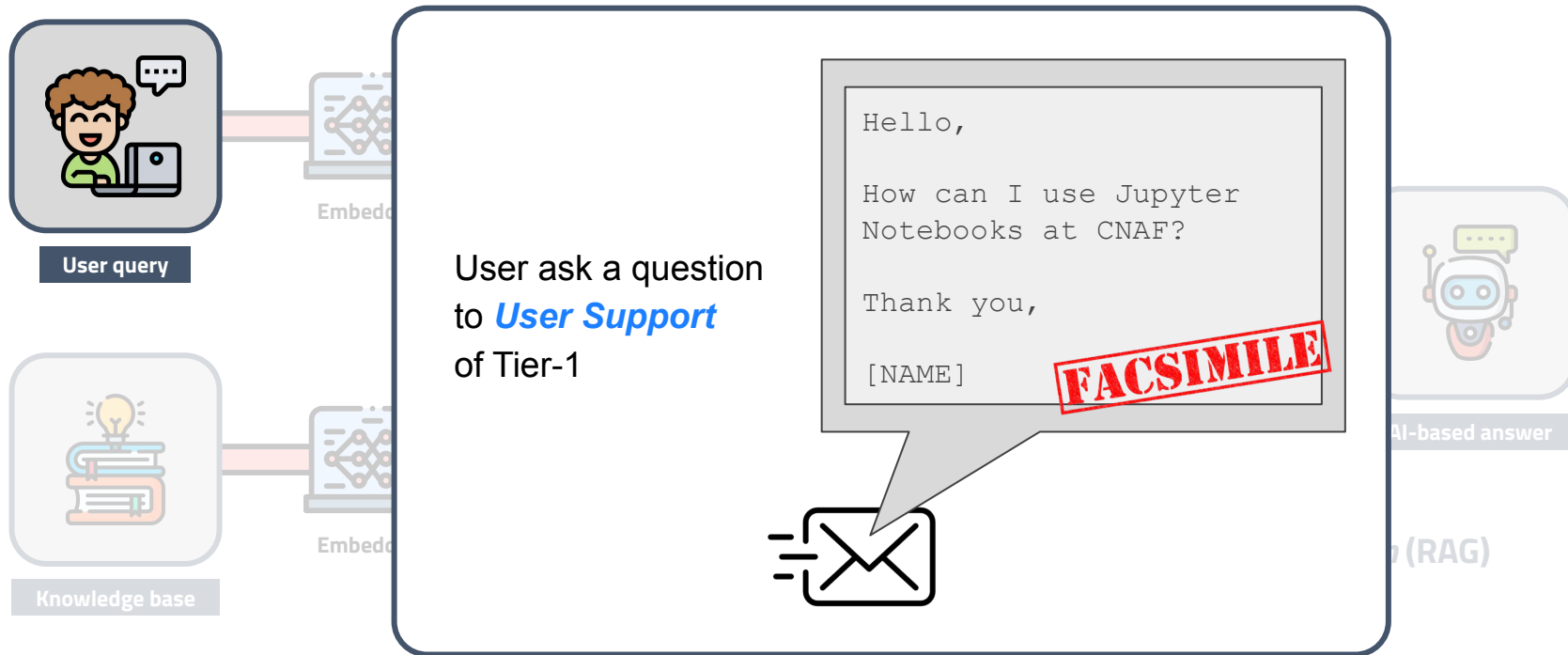
32

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing                    Missione 4 • **Istruzione e Ricerca**

# *How to build a Digital User Supporter*



*Retrieval-Augmented Generation* (RAG)

# How to build a Digital User Supporter

User query

Knowledge base

Embedc

Embedc

User ask a question to *User Support* of Tier-1

Hello,

How can I use Jupyter Notebooks at CNAF?

Thank you,

[NAME]

FACSIMILE

AI-based answer

(RAG)

34

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# *How to build a Digital User Supporter*

The knowledge base (**Tier-1 User Guide**) is mapped to a vector space through an **embedding model**, such as `all-mpnet-base-v2`

**User query**

**Embedding**

**Knowledge base**

**Embedding**

**Vector space**

36

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# *How to build a Digital User Supporter*



**User query**

**Embedding**

**Similarity search**

**Knowledge base**

**Embedding**

**Vector database**

To store the embedding results, a **database** optimized for **handling large vectors** and their **metadata** (e.g., source text) is required. It should also support **efficient searches** within the vector space (features like these are at the core of Chroma)

**Chroma**

# *How to build a Digital User Supporter*



**Embedding**

**User query**

**Knowledge base**

**Embedding**

As well as the knowledge base, also the user request is mapped to a vector space relying on the same **embedding model** (`all-mpnet-base-v2`)

Ciao,

Come faccio a girare un notebook Jupyter al CNAF?

Grazie,

[NAME]

FACSIMILE

**Vector space**

38

# How to build a Digital User Supporter



**User query**

**Embedding**

**Knowledge base**

**Embedding**

**Similarity search**

**Vector database**

To retrieve relevant information from the **guide** to respond to the **user request**, we can compare the embeddings from both cases based on *similarity* (e.g., dot product)



Once the guide embeddings that best match the request embedding are identified, the **vector database** can be queried to **retrieve the corresponding metadata**

39

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# How to build a Digital User Supporter

The retrieved text segments are used to tailor the *template* for the textual **prompt** submitted to the **Language Model**

```
Answer the user's questions
based on the below context as
you were answering to an email
in a professional style.

<context>
{context}
</context>
```

FACSIMILE

**Prompt assembly**

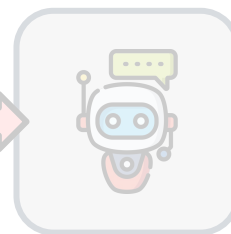Large Language Models

AI-based answer

*Retrieval-Augmented Generation* (RAG)

40

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing          Missione 4 • **Istruzione e Ricerca**

# *How to build a Digital User Supporter*

The generated prompt can then be submitted to a *Large Language Model* (LLM), which, even if not an "expert" in a specific subject, is capable of extracting the **"relevant" content** from `{context}` and rephrasing it

Since fine-tuning is not required, LLMs are often described as **general-purpose** technologies and are thus also referred to as **Foundation Models**

∞ Meta

`ollama run llama3`

Large Language Models

AI-based answer

User que...

Knowledge...

...ly

...val-Augmented Generation (RAG)

41

# How to build a Digital User Supporter



Dear User,

To answer your question directly, yes, you can use Jupyter notebooks through the JupyterHub service at Tier-1. The service is reachable via b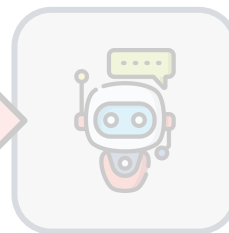rowser at the following page: <https://jupyterhub-t1.cr.cnaf.infn.it/>. Please note that you will need to have an account and be logged in to access the notebooks.

If you have any further questions or concerns, please feel free to ask.

Best regards,
[Your Name]

FACSIMILE

**Note:** automated answer generated with *Llama 2*

AI-based answer

42

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing        Missione 4 • **Istruzione e Ricerca**

# Here are my conclusions!



- **AI is reshaping the computing landscape**, driving technological evolution, influencing market trends, and dominating software development worldwide

- **CNAF should ride the AI wave**, exploring and implementing new technologies to drive its evolution within an AI-dominated world
  - *the development of a Digital User Supporter is a step in this direction, but also **other units can benefit from using AI***

- The results obtained with the first prototype of the **Digital User Supporter** are promising but there is still room for improvements
  - *e.g., multi-expert model, enlarged knowledge base (also including emails)*

- As AI companies develop increasingly large models **requiring massive GPU resources**, it is essential to implement tailored solutions
  - *this is crucial to avoid that AI chatbots **monopolize GPU resources** needed for scientific research*

44

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing        Missione 4 • **Istruzione e Ricerca**

# Thanks!

## Any questions or comments?

**Matteo Barbetti** (INFN CNAF)

*email:* matteo.barbetti@cnaf.infn.it