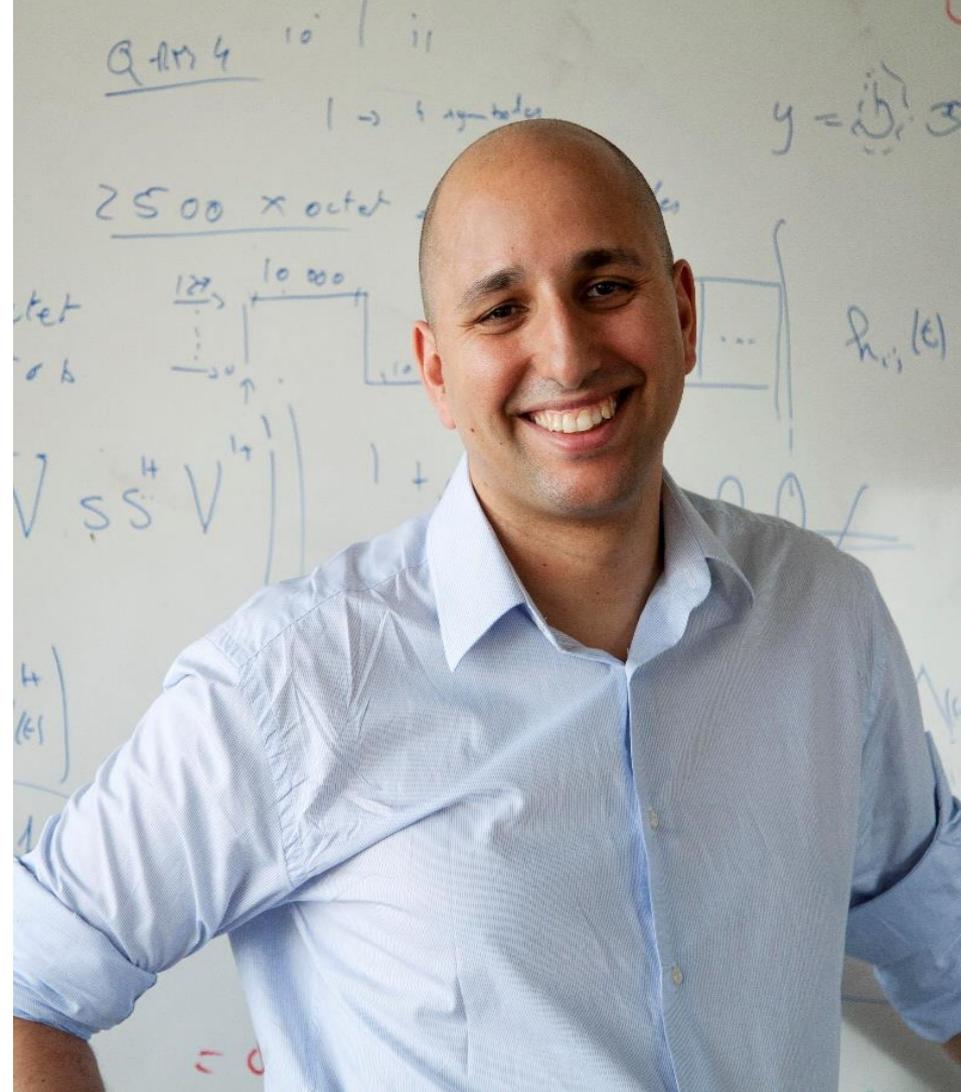# The Super Power of (Open Source) Large Language Models

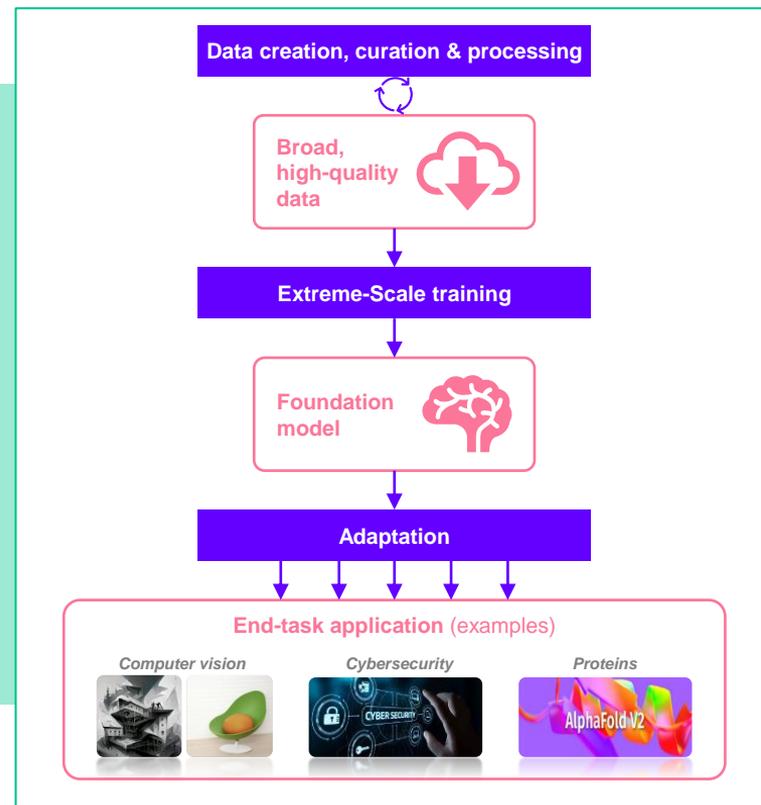## Prof. Merouane Debbah

tii.ae

# About The Researcher

- Chief Researcher at the Technology Innovation Institute

- IEEE, EURASIP and WWRF Fellow

- Citations: 58000+, h-index:108

- More than 35 Best papers Awards

- More than 50 patents

- IEEE Signal Processing Society Distinguished Industry Speaker (2021-2022)

- Field of Research: AI and beyond 5G Systems

# What is LLM and how it works

- LLM stands for **Language-based Learning Machine**

- It is a **type of machine learning model** that uses **natural language** input and output

- LLMs are **trained on large amounts of text data** and **can perform a variety of end-tasks** related to natural language processing (NLP)

- The models use techniques like **neural networks, specifically transformer-based models** like BERT, GPT, T5... which have proven to be very effective in natural language tasks



**Data creation, curation & processing**

Broad, high-quality data

**Extreme-Scale training**

Foundation model

**Adaptation**

**End-task application** (examples)

*Computer vision*   *Cybersecurity*   *Proteins*

AlphaFold V2

# Large Language Models (LLMs) are eating machine learning

## LLMs provide a universal text-based interface to tackle any tasks:

Text instructions → **LLM** → Answer (text output)

What companies were mentioned in this article?
Write an essay based on these notes:
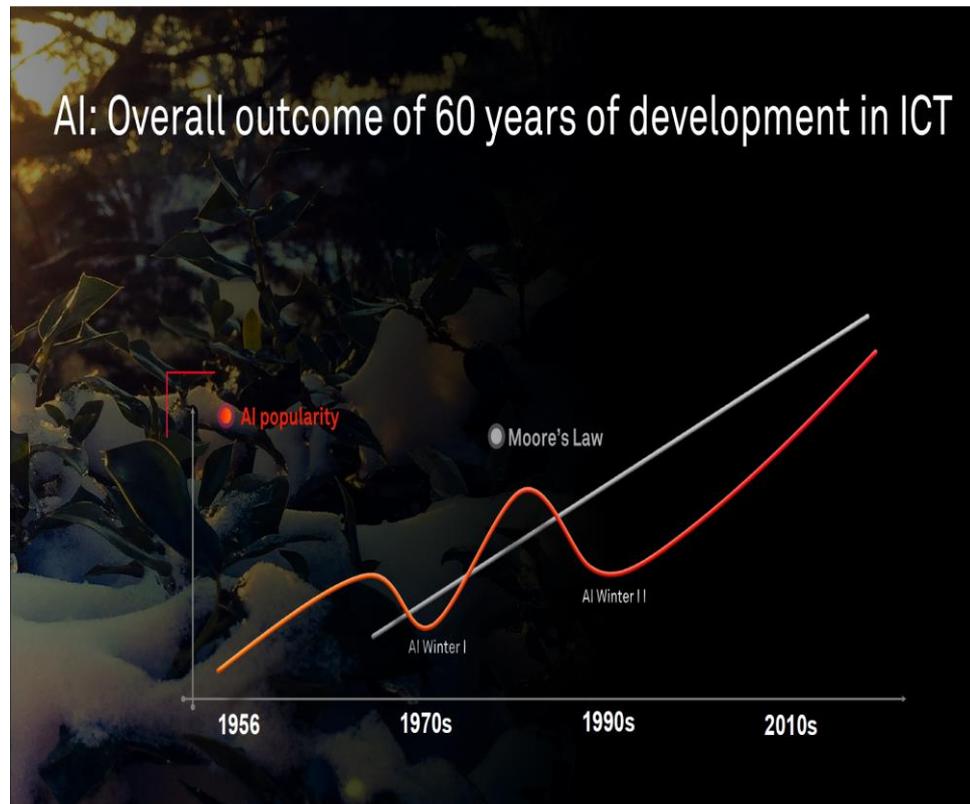Write an Instagram ad for...

e.g. GPT-3

**Key aspects of LLMs:**

🧠 They are generalists, able to tackle broad tasks just from instructions.

📈 Their capabilities increase as you scale-up in size/compute.

🧪 One of the main business & research interest in machine learning.
from Google, DeepMind, Microsoft, etc. + large start-ups such as OpenAI and Cohere.

# Why now?

**Massive amounts of data** that can be used to train Machine Learning models are being generated, for example through daily creation of billions of images, online click streams, voice and video, mobile locations, and sensors embedded in the Internet of Things devices.

**Computing capacity** has become available to train larger and more complex models much faster. Graphics processing units (GPUs), originally designed to render the computer graphics in video games, have been repurposed to execute the data and algorithm crunching required for machine learning at speeds many times faster than traditional processor chips.

**Machine-learning algorithms** have progressed in recent years, especially through the development of deep learning and reinforcement-learning techniques based on neural networks.



AI: Overall outcome of 60 years of development in ICT

# What LLM can do

## Content creation

**Text / Code**
- Idea generation
- Text writing (book, training, course, plan…)
- Copywriting (emails, ads, blog posts…)
- Code writing (website, app…)

**Visuals / Sounds**
- Image generation (text-to-image)
- Video generation (text-to-video)
- Voice generation (text-to-voice)
- Game design (AR, 3D design…)

## Content curation and analysis

**Rewrite / Summarize**
- Text summary
- Video summary
- Audio transcript
- Language translation
- Information clustering / formatting

**Extract**
- Information retrieval
- Web search / Benchmarking
- Q&A

**Analyze**
- Data analytics and forecast
- Visual analytics
- Sentiment / Intent recognition (ex. Fraud)

## Task automation

- Chatbot / Virtual assistant
- Scheduling (meetings, tasks…)
- Text editing (spell check, paraphrasing)
- Visuals editing (video cutting, image editing…)
- Data cleansing
- Code auditing
- Robot control

# Noor Released in April 2022



## Technology Innovation Institute Announces Launch of NOOR, the World's Largest Arabic NLP Model

13 Apr, 2022

Model is most powerful one in Arabic language to date with 10 billion parameters

Features applications in automated summarization, chatbots, personalized marketing

# Noor
## Released April, 2022

*Arabic language model with 10B parameters*

## First published
## Arabic Article by Noor
## November, 2022

*Noor Arabic language model generates an article*

# Noor - Narrator

Narrator

العودة إلى الصفحة الرئيسة

راوي

الرجاء إضافة نص يتكون من 5 حرف على الأقل

كرة القدم رياضة شعبية

21 حرف

انشاء نص

النص المنتج

ترقب التلخيص هنا

Made with ♥ At TII

# Noor - Summarization

Back

# Falcon meets LLama

# Abu Dhabi makes its Falcon 40B AI model open source

By Lisa Barrington ⌄

May 25, 2023 12:02 PM GMT+4 · Updated 4 days ago

DUBAI, May 25 (Reuters) - The emirate of Abu Dhabi is making a large-scale artificial intelligence model, "Falcon 40B", available open source for research and commercial use, the government's Advanced Technology Research Council (ATRC) said on Thursday.

ATRC's commercial investment arm VentureOne said it would also back viable ideas that come from using the model.

Falcon 40B is a foundational large language model (LLM) with 40 billion parameters and trained on one trillion tokens which was developed by the Technology Innovation Institute (TII), a research centre within ATRC.

## Posts by Thomas

**Thomas Wolf** · 1st
Co-founder at 🤗 Hugging Face
2d · 🌐

LLaMa is getting dethroned 👑 There is a fresh new pretrained LLM sitting on the top of the Open LLM Leaderboard: Falcon 40B 🚀 from the Technology Innovation Institute

And it has some special features:
- an architecture strongly optimized for inference (MQA, flash-attention, parallel attention/MLP, smaller than 65B)
- open-source with a special licence allowing commercial use (an interesting approach check it out on the model page - interesting to see the reception)
- available in two sizes: 40B and 7B parameters

More details here:
- Falcon Model 40B: https://lnkd.in/efcT7-9q
- Falcon Model 7B: https://lnkd.in/eBuCyEgq
- Technology Innovation Institute HF org page: https://lnkd.in/e9Vfc5ST
- The Open LLM Leaderboard: https://lnkd.in/eUzaDuUu

Super interested to see what the community will think about this new model!
Don't hesitate to comment below 👇

UAE's Falcon 40B Dominates Leaderboard: Ranks #1 Globally in Latest Hugging Face Independent Verification of Open-source AI Models



TII

Falcon 40B
Ranked #1 Globally
on Hugging Face  Open LLM Leaderboard

## Falcon 40B

# Falcon LLM

- Foundational large language model (LLM) with 40 billion parameters trained on one trillion tokens.

- Utilizes only 75% of GPT-3's training compute, outperforming it with a fifth of the inference compute.

- Custom tooling and unique data pipeline for high-quality content extraction from the web.

- Optimized architecture for performance and efficiency, matching state-of-the-art LLMs.

- Pretrained on a dataset of five trillion tokens, including web crawls, research papers, and social media conversations.

- Applications include chatbots, customer service, virtual assistants, translation, content generation, and sentiment analysis.

- Aims to automate repetitive tasks, enhancing efficiency for Emirati companies and individuals' daily lives.

# Introducing Falcon LLM

Technology Innovation Institute has open-sourced Falcon LLM for research and commercial utilization.

**Access Falcon LLM** →

TII is calling for proposals from the global research community and SME entrepreneurs to submit use cases for Falcon LLM.

**Submit Use Case Proposal** →

# LLM market landscape

**A limited number of major players as cornerstone platforms…**

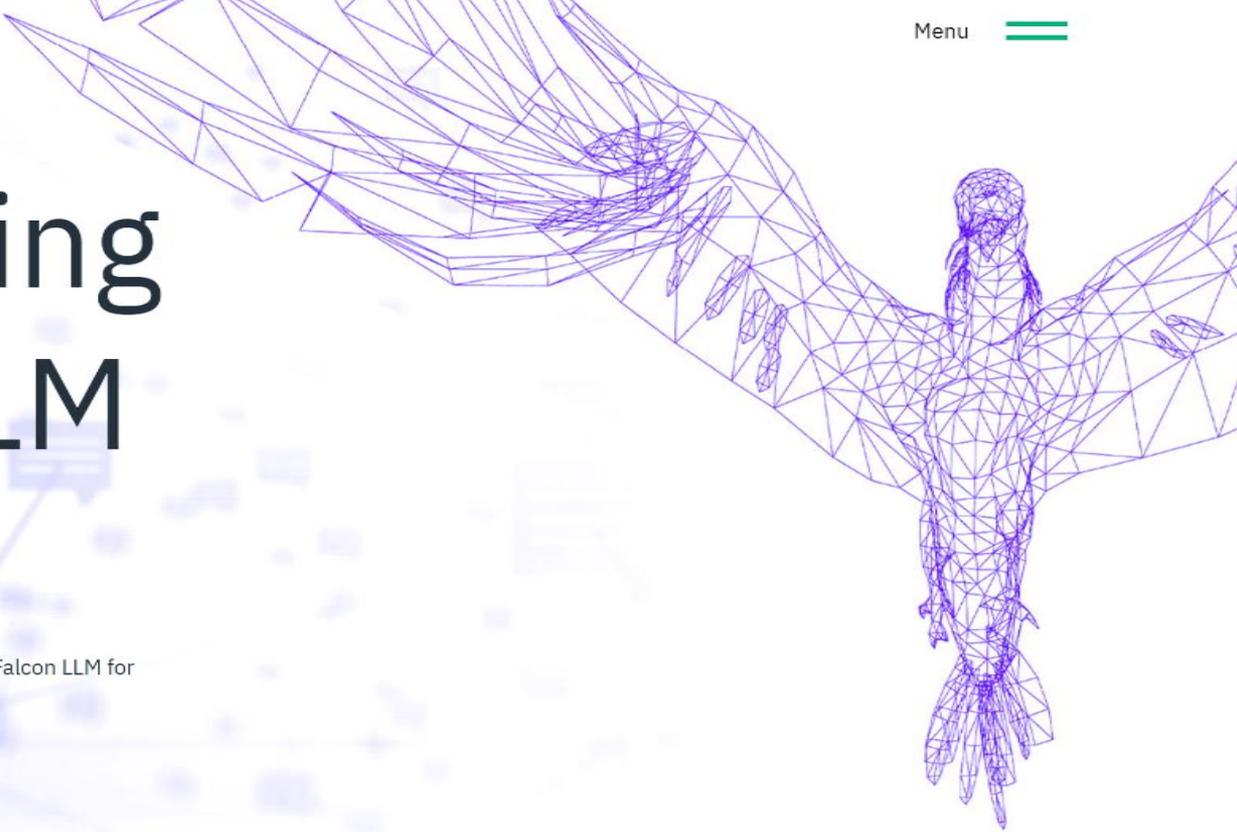| Platform | G Google | OpenAI | AI21 labs | DeepMind | NVIDIA | + others |
|---|---|---|---|---|---|---|
| **Model** | LaMDA | GPT-3 | Jurrasic | Gopher | Megatron-Turing | Inflection, Meta, cohere, Tencent, Baidu, EleutherAI |
| **SIZE** [Bn parameters] | 137 | 175 | 178 | 280 | 530 | |
| **Training tokens** [bn] | 168 | 300 | 300 | 300 | 270 | |

**… Supporting hundreds of startups focused on end-user applications**



*Grid of startup logos organized by category:*

**Text:** Smartwriter.ai, Hypertype, Lately, Autobound, Writesonic, Jasper, cogram, genei, YOU, copysmith, AI21labs, letterdrop, Creatext, jenni, mavenoid, anyword, PERSADO, frase, regie.ai, Mutiny, Linguix, Hypotenuse AI, WRITER, OTHERSIDEAI, copy.ai, copymatic, COMPOSE AI

**Image:** ClipDrop, p-e-n-cil, beautiful.ai, PhotoRoom, BRIA, Facet, Poly, CSM, Blend, HYPAR, ROSEBUD.AI, maket, KAEDIM, BOTIKA, Autoenhance.ai, Sloyd, MODULIZE, Re:cast AI, uizard, Imagen, Hexo AI

**Audio:** MURF.AI, REPLICA, notably, Endel, WELLSAID, AssemblyAI, DEEPGRAM, krisp, Speechify, RESEMBLE.AI, MIMI, KAIZAN, coqui, Mubert, Neural Space, soundful, PODCASTLE, moises, VOICEMOD, Listnr, LOVO, Vocal Clarity, Dubverse, AD AURIS

**Video:** ZUBTITLE, TERRA, Peech, VOCHI, Maverick, recast, VEED.IO, Basch.io, inworld, deepdub.ai, WOMBO, tavus, FATHOM, runway, Maria, XEMBLY, PICTORY, vidyo.ai, Rephrase.ai, lumen5, Steve AI, windsor.io, YEPIC, Colossyan, METAPHYSIC, Potion

**Code:** Debuild, tabnine, Codiga, Locofy, AIXcoder, Mintlify, maya, MutableAI, Cod.i, durable, The.com, bloop, replit, ENZYME, codota, DhiWise, CODACY, anima, warp, Metabob

**Chatbots:** lang.ai, PolyAI, Tymely, Incentivai, Kasisto, ushur, MLY, CRESTA, Elise.AI, verloop.io, Replika, ultimate.ai, Cohere, Sapling, haptik, ada, Forethought, OBSERVE.AI, XOKind, Balto, Certainly.

**Search:** glean, Looria, Hebbia, consensus, Air, vectara, Pinecone, qdrant

**Gaming:** charisma.ai, hidden door, LATITUDE, Spellbrush

**Data:** Pilot, gretel, DATAHERALD, SYNTEGRA, Mirry.AI, BIFROST, datagen

**ML platforms:** slai, Jina, symbl.ai, Adept, aporia, GANTRY, deepset, Synthesis.ai, Archistar, Galileo, featureform

# Key considerations when building LLMs

**Data**

The more data the model is exposed to (volume and diversity), the better it will perform on unseen examples

**Model architecture**

Choice of model architecture is crucial for the LLM's performance. Transformer-based models like BERT, GPT & T5 have proven to be very effective in natural language tasks

**Pre-processing**

Cleaning and preparing data which includes tasks such as tokenization, stemming, and lemmatization

**Fine-tuning**

To adapt for specific tasks, it is often necessary to add task-specific layers to the pre-trained model

**Evaluation**

Careful assessment of the model to ensure that it is performing well on the task it was designed for

**Deployment**

Different deployment options can be considered depending on the use-case, such as cloud-based services, on-premises or edge devices

**Explain-ability**

Models can be complex, and it can be hard to understand how they arrived at their predictions. There are methods like LIME & SHAP to make the predictions interpretable

**Privacy & Security**

Since LLM models may handle sensitive information, it is important to consider privacy and security issues when building and deploying them

**Bias & Toxicity**

LLM models can inadvertently perpetuate and amplify societal biases (ex. People minorities representation) and toxicity (ex. hate speech or offensive language) present in the training data

**Mitigation techniques**

- Use representative and diverse training data
- Carefully evaluate and test the model
- Use techniques like debiasing & Fairness-Aware learning
- Monitoring the model's output in production

# Thank You

tii.ae