MULTIVERSE
COMPUTING

AISC

10

JAVIER
ALONSO
MENCÍA

ML ENGINEER
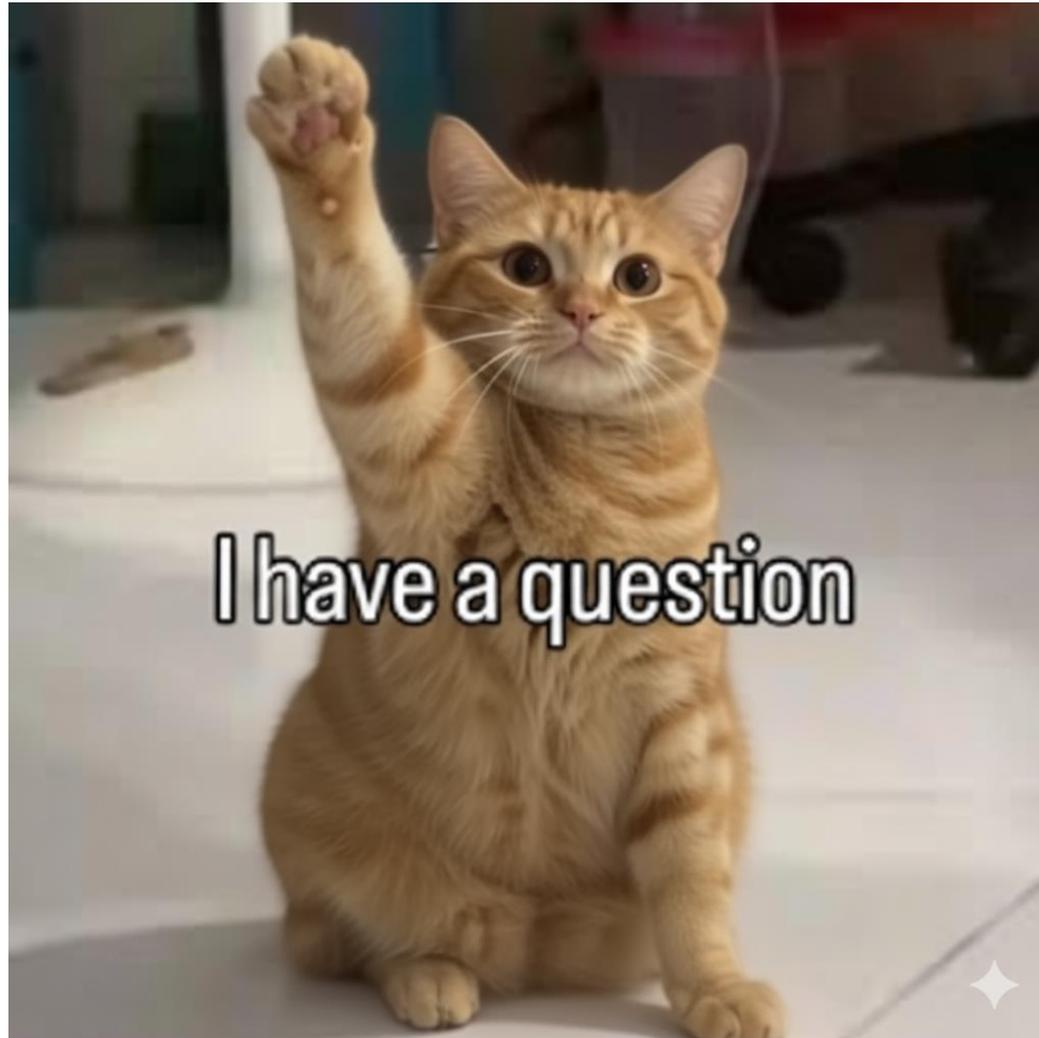
18

**Machine Learning en el mundo real:**

lo que no te cuentan en clase

👋 Quién soy

🚀 Mi trayectoria

🧠 Qué hacemos en MC

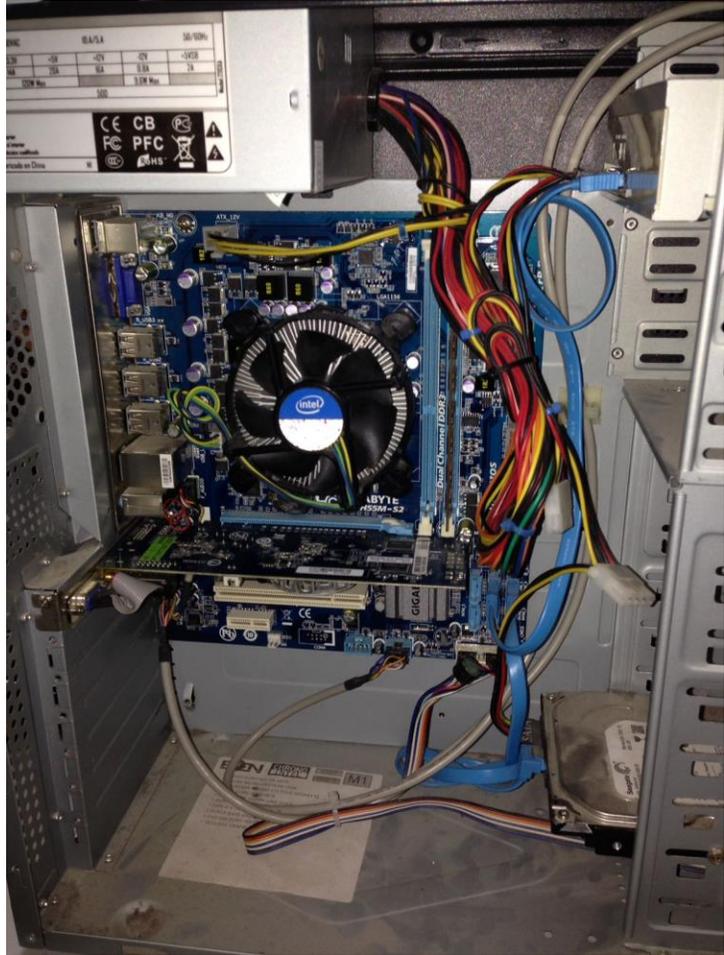🤖 Vivir y trabajar en ML

🌟 Consejos para tu futuro

**Javier Alonso Mencía**
Machine Learning Engineer
Multiverse Computing

# Ir a clase es solo el comienzo

Y ahora qué?

# ADELA: a conversational virtual assistant to prevent delirium in hospitalized older persons

Javier Alonso-Mencía[1,2] · Marta Castro-Rodríguez[3] · Beatriz Herrero-Pinilla[3] ·
Juan M. Alonso-Weber[1] · Leocadio Rodríguez-Mañas[3] ·
Rodrigo Pérez-Rodríguez[4]

## Abstract

Delirium is a sudden mental state that causes confusion and disorientation, affecting a person's ability to think and remember clearly. Virtual assistants are a promising alternative for non-pharmacological interventions. This research aims to present a prototype of ADELA, a conversational assistant to prevent delirium in hospitalized older persons who speak Spanish. A co-creation process with medical experts to

11

1/56

# Piso en venta en Almagro

Chamberí, Madrid  ⦿ Ver mapa

**3.780.000 €**

233 m²  │  3 hab.  │  Planta 4ª exterior con ascensor

Lujo

♡ Guardar        🗑 Descartar        ↱ Compartir

# Time Line



| 2020 | 2021 | 2022 | 2023 | 2024 | 2025 |
|------|------|------|------|------|------|
| uc3m | uc3m | EA | idealista | idealista | MULTIVERSE COMPUTING |
| Bachelor | FUNDACIÓN DE INVESTIGACIÓN BIOMÉDICA SaludMadrid Hospital Universitario de Getafe | Data Science intern | Data Scientist | Data Scientist | Machine Learning Engineer |
|  | Masters Researcher |  |  |  |  |

# Key Highlights

## TEAM

**+180** people

**+35** nationalities

**20%** PhD[1]

**30%** women

## FINANCIAL

**$65M+**
Min 2025 AnRR
Targeting $**150M+ by EOY**

**17x**
**YoY AnRR**
Growth

## IP

**+200** Patents[2]

**+40** Research Publications

## FUNDING

**$250M**
Series B closed June 2025,
**the largest quantum-AI round in Europe**

# Our locations

**Spain - HQ**

Canada

USA

UK

Germany

France

# Application Sectors



**Finance**

**Manufacturing**

**Energy**

**Healthcare & Life Sciences**

**Chemistry**

**Cybersecurity**

**Hydrogen**

**Defense**

**Pure Engineering**

**Aerospace**

**Others**

Todo esto está muy bien... pero, ¿qué hacéis en Multiverse?

# AI Model Compression

Harness the Power of Tensor Networks

**CompactifAI**

AI Model Compressor

**Reduce the number of parameters** of AI models without compromising accuracy, making them more accessible, affordable and sustainable.

**+**

**Cost-saving**

Slash Computational and Infrastructure Costs

**Private**

Run Anywhere: on Premise, Cloud, or Any Device

**Small**

Less GPU Memory and Storage

**Fast**

Faster training and Inference

**Efficient**

Less Energy Consumption and Less $CO_2$ emission

# Skyrocketing AI Compute Costs

The cost of training notable AI models has grown by a factor of 2.4x per year for the past eight years, **suggesting that the largest models will cost over a billion dollars by 2027.**

**Amortized hardware and energy cost to train notable AI models over time**

Ben Cottier, Robi Rahman, Loredana Fattori, Nestor Maslej and David Owen (2024), "How Much Does It Cost to Train Frontier AI Models?". Published online at epoch.ai. Retrieved from: 'https://epoch.ai/blog/how-much-does-it-cost-to-train-frontier-ai-models' [online resource]. Updated on January 13th , 2025

# MULTIVERSE COMPUTING

Multiverse Computing is a **global leader in compressed AI** powered by quantum-inspired tensor networks, universalizing affordable and efficient AI across cloud, edge or on-prem.

**TOP100 Fastest-growing startups in Europe**

\sifted/

2025

**'Future Unicorn' Award**

DIGITALEUROPE

2024

---

**We developed CompactifAI, an AI model compression technology** to enhance AI system performance by reducing LLM size

**CompactifAI**
Multiverse Computing

**98%** compression rates of LLMs[1]
**4x-12x** speed
**50-80%** cost savings

For some clients, this translates to **half-a-billion $ savings p.a.**

---

We offer the market's **leading AI models as**

**LLaMA** by ∞Meta   **OpenAI**   **deepseek**

**MISTRAL AI_**   **ultralytics YOLO**

and more...

**+ we launched our CompactifAI API in 2025,** offering original and compressed models at an **unbeatable price**

**Available on**

aws marketplace

---

We count with leading in-house compressed **NanoModels to run anywhere**

**SuperFly**   **94M** parameters

**ChickBrain**   **3.2B** parameters

**Deliver top-tier performance** on ultra-light hardware, enabling **true on-device and edge AI**

**Buzzy AI startup Multiverse creates two of the smallest high-performing models ever"**

TechCrunch

# Welcome to the CompactifAI API

Build intelligent applications with our state-of-the-art language models

Quickstart →        API Reference 📖



CompactifAI
AI Model Compressor

MULTIVERSE
COMPUTING

# Get started with CompactifAI

### Introduction

Learn about the CompactifAI API and its capabilities

### Pricing

See our pricing for the CompactifAI API

pero.. ¿cómo hacéis la compresión?

# Compression Steps

**Original Model**

**Compressed Model**

**Profiling**

Investigating the Model
+
Analyze Sensitivity of Layers
to Compression

**Compression**

Tensor Network Compression

**Healing**

Small Training to Recover Accuracy

# Speed

**From the first token to the full response — in record time**. CompactifAI significantly reduces latency and boosts throughput, delivering faster performance for real-time applications, even on resource-constrained hardware.

## Llama 3.3 70B
### Time to First Token[1]



## Microsoft Phi 4
### Time to First Token[1]



## Llama 3.3 70B
### Tokens per Second[1]



## Microsoft Phi 4
### Tokens per Second[1]



↓**53%**
**Reduction in TTFT**
vs. original model

↓**61%**
**Reduction in TTFT**
vs. original model

Up to
**2.4X**
**Increase in Speed**
vs. original models

Llama 3.3 70B compressed model by MVC offers a **nearly 145% faster throughput** than the original Meta LLM.

[1] Real-time performance evaluation for MVC Compressed vs Original models. Source: External benchmarks by a leading enterprise hardware provider, 2025.

# Cost Efficiency & Sustainability

**Achieve more with less**. Cut your AI expenses and lower your carbon footprint using ultra-efficient LLMs. Our compressed models operate on smaller hardware, consume less power, and provide enterprise-level performance.

## Llama 3.1 8B

### Energy Consumption[1]

Energy (kWh) at stage: prefill

- Llama 3.1 8B Original
- Llama 3.1 8B Slim by CompactifAI

Y-axis: 0.0E+00, 5.0E-07, 1.0E-06, 1.5E-06, 2.0E-06, 2.5E-06, 3.0E-06

X-axis: CPU, RAM, GPU

## Llama 3.3 70B

### RAM Usage[2]

MB — Y-axis: 0, 20, 40, 60, 80, 100, 120, 140, 160, 180

X-axis: Original, Multiverse

CompactifAI

## Microsoft Phi 4

### RAM Usage[2]

MB — Y-axis: 0, 5, 10, 15, 20, 25, 30, 35

X-axis: Original, Multiverse

CompactifAI

↓**50%**
**GPU energy consump.**
vs. original model

CompactifAI's slim model **halves the energy demand of the original Llama 3.1 8B**, making it ideal for sustainable, cost-effective AI deployments.

Up to
↓**70%**
**Average RAM usage**
vs. original models

**MVC models deliver significantly better speed and comparable accuracy** — achieving faster inference and stronger results without compromising performance.

# Accuracy

**Exceptional performance guaranteed**. Our Slim models consistently match or outperform the accuracy of original models on key benchmarks—demonstrating that smaller truly means smarter.

## Llama 3.1 8B



Legend: Llama 3.1 8B | Llama 3.2 3B | Llama 3.1 8B Slim by CompactifAI

| Benchmark | Llama 3.1 8B | Llama 3.2 3B | Llama 3.1 8B Slim |
|---|---|---|---|
| MMLU Pro | 37 | 32 | 41 |
| MATH-500 | 29 | 15 | 76 |
| GSM8K | 76 | 65 | 77 |
| GPQA Diamond | 30 | 21 | 31 |

## Llama 3.3 70B



Accuracy[2]: Original ≈65%, Multiverse ≈62%

*CompactifAI*

## Microsoft Phi 4



Accuracy[2]: Original 70%, Multiverse ≈61%

*CompactifAI*

Up to
↑ **107%**
**Higher Accuracy**
vs. Meta's model

CompactifAI Slim model outperform Llama 3.2 3B by **more than double in accuracy** on key reasoning tasks like MATH-500.

On average
**6%**
**Minimal Accuracy Drop**
vs. Original models

On average, MVC models maintain comparable accuracy to original versions, with only a minimal drop of up to 6 percentage points — while **delivering superior performance in speed.**

[1] Source: MVC internal benchmark, 2025
[2] Source: External benchmarks by a leading enterprise hardware provider, 2025.

28

¿Solo comprimis modelos?
¿Quién los usa?

# Success Stories

# Compressed LLM Application in Telecommunications

Use case: Internal Customer Service Chatbot for 8000 Sales Representatives

**Llama 3.1 8B**

**Industry:** Telecommunications

**Client:** IBEX 35 telecommunications provider with global presence

**Goal:** Compress LLM to use it within an internal RAG chatbot for agents in client's stores.

| Cost reduction | Speed |
|---|---|
| ↓50% | ↑46% |
| incl storage requirements cut by 78%, energy cost cut by 63%, inference cost down 29% | Reduced time to first token, plus extra 26% reduced in-token latency |

| Compression achieved | Accuracy loss |
|---|---|
| 80% | 0% |
| 60% fewer parameters plus quantization | Vs original model when using optimized prompts |

**Reduced costs**

associated with the use of expensive API services (GPT, Gemini, etc.)

**Speedup**

Reduced latency – faster responses

**Efficiency**

Low power consumption

**Security**

Increased data privacy

# Compressed Computer Vision Application in Defense & Military (1/2)

Use case: Satellite Image Processing of 670K km2 for Object Detection

**YOLO v8-x**

**Industry:** Defense

**Client:** Transnational defense and security organization

**Goal:**

Accelerate inference and reduce infrastructure costs.

Keep model accuracy as high as possible.

Compression of the computer vision model **YOLOv8-x.**

Use it to process high-resolution satellite imagery.

- 10 cm/pixel resolution.

- 670,000 km2 surface to analyze.

- Refresh rate: 4 times per hour.

Images in 3 different spectral bands: RGB, Infrared, SAR

Tensor network compressibility of convolutional models
https://arxiv.org/pdf/2403.14379

# Compressed Computer Vision Application in Automotive

Use case: On-Edge In-Car Virtual Assistant Using Compressed TextToSpeech AI Model

**Style TTS**

**Industry:** Automotive

**Client:** Leading European automotive manufacturer

**Goal:**

- Compress the model as much as possible with no noticeable loss in the quality of the assistant's voice.

- Preliminary results looking promising -37% parameter reduction with no drop in audio quality.

**Next Steps**

Continue compressing further, and combine CompactifAI with other techniques such as quantization.

| | Base Model<br>Style TTS | CompactifAI<br>Multiverse Computing<br>Style TTS |
|---|---|---|
| **Parameter** Reduction | - | ↓37.25% |
| **Memory** Reduction | - | ↓29.05% |
| **Accuracy:** NISQA MOS | 4.507 | 4.553 |
| **Accuracy:** NISQA Noise | 4.275 | 4.134 |
| **Accuracy:** DNSMOS | 4.017 | 4.019 |
| **Audio Sample** | 🔊 | 🔊 |

pero... ¿en qué trabajas tú?

# Ser ML Engineer en Multiverse Computing - cómo trabajamos

Ideación y diseño del proyecto

↓

Preparación de datos

↓

Compresión y reentrenamiento

↓

Evaluación con benchmarks

↓

Deployment y soporte al cliente

Compresión de modelos

Proyectos con clientes

# Equipo, herramientas y retos en Multiverse Computing

👥 **Equipo y colaboración**
→ Estructura multidisciplinar: MLOps, DevOps, Research, Servicios.
→ Comunicación constante entre equipos.

🧰 **Herramientas y lenguajes**
→ Python, PyTorch, frameworks de NVIDIA.
→ LangChain, MindEx y librerías de agentes.
→ Infraestructura y MLOps internos.

⚡ **Retos y habilidades**
→ Técnicos: escalabilidad, eficiencia, coordinación.
→ Personales: adaptación, comunicación, gestión del tiempo.
→ Claves: curiosidad, pensamiento crítico, aprendizaje continuo.

¿Es un buen momento para estudiar ML/IA?

# Retos del sector tecnológico

⚙️ **Retos actuales**

→ **Boom de la IA:** cientos de startups, inversión y crecimiento acelerado.
→ **Incertidumbre:** el auge sigue, pero nadie sabe cuánto durará (**burbuja ¿?)**
→ **Avance tecnológico rápido**, aunque el paradigma se mantiene estable.
→ **Desafíos éticos y de talento:** falta de profesionales y necesidad de una IA responsable.

💬 **Lo que me ha sorprendido/cambiado mi forma de pensar**

→ Las **habilidades sociales** son tan importantes como las técnicas.
→ **Adaptarse a cada persona y contexto** es clave.
→ La colaboración y la empatía marcan la diferencia en los proyectos.

# Cómo mantenerse actualizado

🔁 **Aprender siempre**
La tecnología cambia cada pocos años: seguir aprendiendo es parte del trabajo.

⚙️ **Actualizarse constantemente**
Surgen nuevos modelos, herramientas y enfoques continuamente.

🧪 **Probar cosas nuevas**
Participar en proyectos, cursos o retos técnicos mantiene tu mente activa.

🤝 **Aprender de otros**
Eventos, comunidades y networking te mantienen al día y conectado.

¿Y si volviese a empezar?

# Qué haría diferente y habilidades clave

## 🕰️ Qué haría diferente

→ Repetiría una carrera similar (**IA / Datos / Informática / etc**)

→ Haría un **máster** y luego consideraría un **doctorado fuera de España** por mejores condiciones.

## 🌱 Habilidades clave (soft skills)

→ **Comunicación y habilidades sociales:** ir a charlas, hablar con ponentes y participar en eventos.

→ **Networking:** crear contactos y aprovechar oportunidades.

→ **Proyectos personales:** diferénciate y demuestra iniciativa.

→ **No todo es técnica:** las soft skills abren más puertas que un currículum impecable.

# Máster, prácticas y primer trabajo

🎓 **Elige con propósito**
El sector es competitivo, elige un máster o proyecto que te **haga crecer**.

💼 **Empieza con prácticas**
Las prácticas son la mejor puerta de entrada al sector.

🧭 **Gana experiencia real**
Una vez dentro, es más fácil conseguir un puesto fijo o cambiar de empresa.

🚀 **Haz lo que te motive**
No elijas por nota: busca lo que **te apasione y te rete**.

🧪 **Aprende haciendo**
Crea mini proyectos, investiga y prueba nuevas herramientas.

# One last advice

In a world driven by AI, your greatest advantage is still your ability to connect, communicate, and stay human.

*En un mundo impulsado por la IA, tu mayor ventaja sigue siendo tu capacidad para conectar, comunicarte y mantener tu humanidad.*

*"La curiosidad es la clave del éxito"*
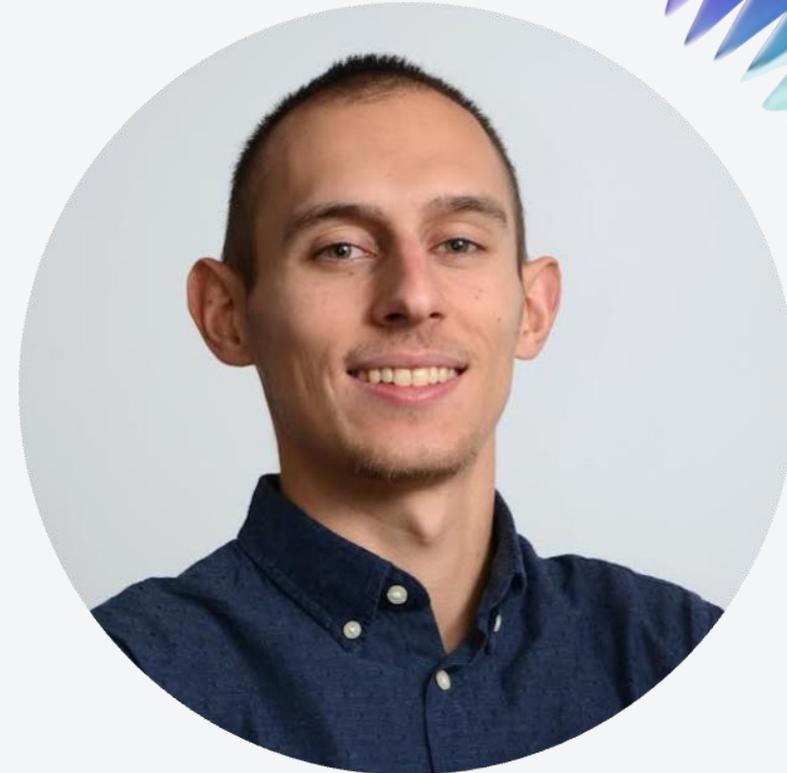
**MULTIVERSE**
COMPUTING

https://multiversecomputing.com/join-us

@javieralonso9

javilonso9@gmail.com

**Javier Alonso Mencía**
Machine Learning Engineer