

# #CONF'

*Barbara DELACROIX  
& Marvin SANT*

*Fondateurs @ Devana.ai*



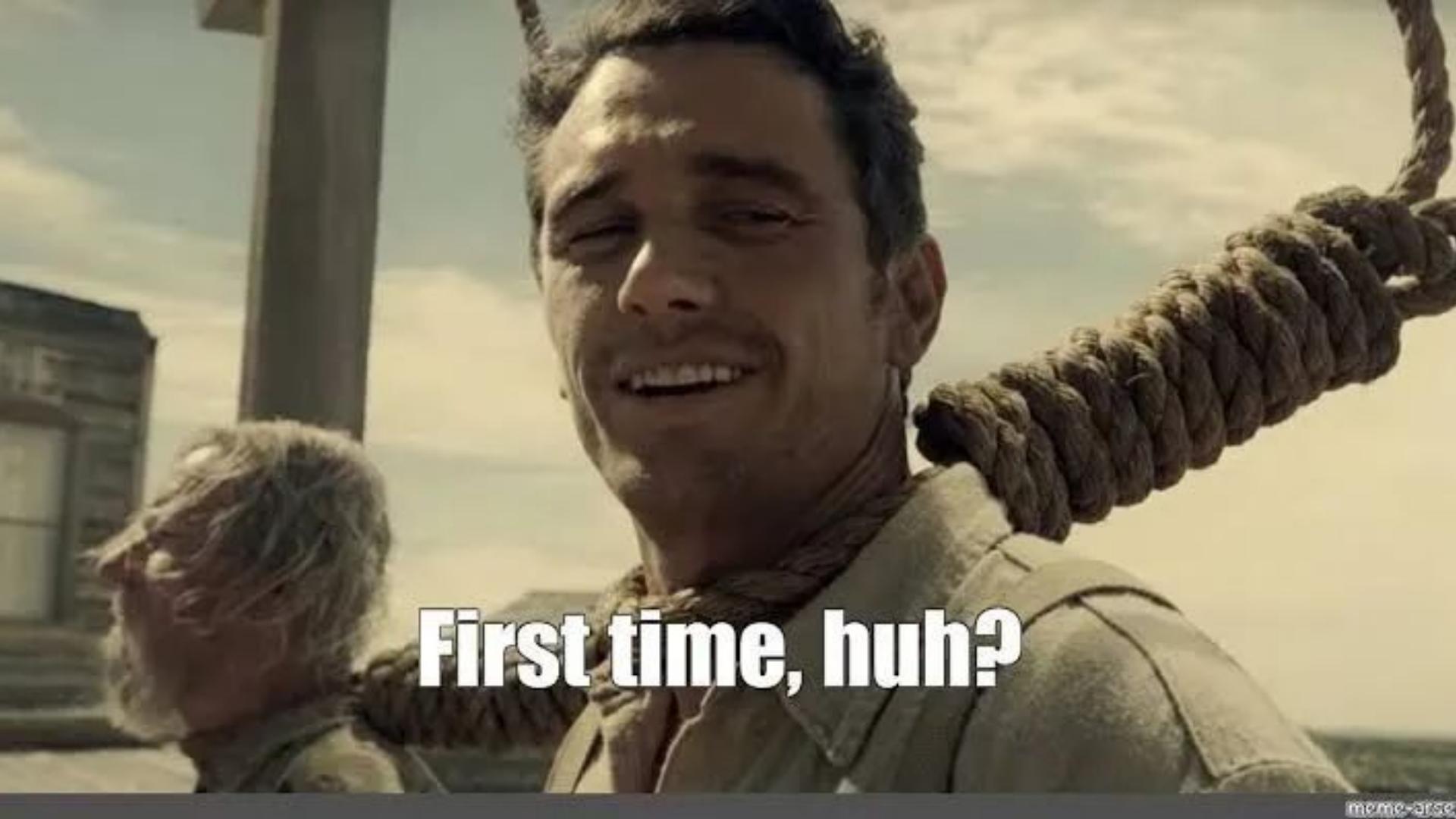
*Le 19/12/23 à 19h*



*le wagon*

*"Le RAG : booster de connaissances pour la GenIA ?"*



A close-up photograph of a man with dark hair and a beard, smiling broadly. He is wearing a light-colored, collared shirt. A thick, dark rope is wrapped around his neck, with one end resting against his shoulder and the other extending towards the top right corner of the frame. In the bottom left corner, the back of another person's head and shoulders are visible, showing short, wavy hair and a greenish-brown jacket.

**First time, huh?**



# Schedule



1- News (15min)



2- Talk (1h30)



3- Enjoy ;)

# Sponsor turbo-platinum



**le wagon**



**Accélérez votre  
carrière.**

**Formez-vous aux  
métiers de la tech.**



# Notre vision

**Nous croyons que chacun(e)  
peut s'épanouir continuellement  
dans sa carrière.**



## Reconversion

Bootcamps intensifs en Web & Data  
400h



## Montée en compétences

Skill Courses à temps partiel  
40h

# À propos du Wagon

Le Wagon est un des leaders mondiaux des formations immersives tech depuis 2013.

- ✓ Débutez une nouvelle carrière
- ✓ Obtenez des compétences techniques
- ✓ Cours en ligne ou sur campus
- ✓ Formats temps plein et temps partiel
- ✓ Formations Data, Développement Web & No-code

**23,000+**

Diplômés dans le monde

**#1**

La formation la plus reconnue

**4,98/5**

Sur +6,600 avis

**40+**

Villes sur tous les continents

**200+**

Startups fondées par nos alumni

**\$1B**

Levé par nos alumni

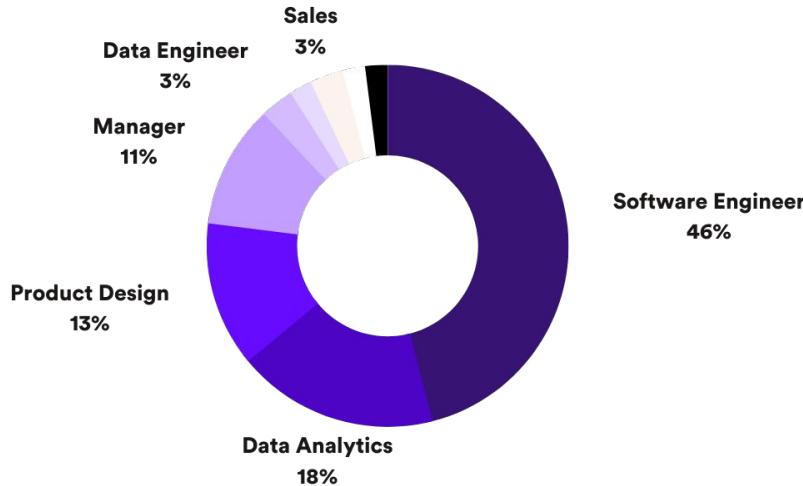
# Présence mondiale



 Campus du Wagon

# L'employabilité après nos formations

Les données sont plus parlantes que les mots.



le wagon



**34 jours**

En moyenne pour décrocher un emploi dans le secteur de la Tech en Europe.



**93%**

De nos diplômés des campus Européens ont trouvé un job dans une entreprise Tech, se sont lancés en Freelance ou ont lancé leur propre startup.



**45k€**

Salaire médian en Europe pour leur premier emploi post-bootcamp.

# Nos bootcamps

200-400h

Pour accéder à de nouvelles opportunités grâce aux formations immersives en développement web et en data.

✓ 1 professeur pour 7 étudiants

✓ Temps partiel ou Temps plein (de 2 à 6 mois)

✓ Sur campus ou en ligne



## Développement Web

Devenez développeur.se et créez des applications web, de la base de données à l'interface utilisateur.



## Data Science & IA

Obtenez les compétences fondamentales du Data Scientist et créez vos propres modèles d'IA.



## Data Analytics

Apprenez à transformer vos données en informations exploitables et décrochez un poste de Data Analyst.



## Data Engineering

Concevez des pipelines de données et développez des bases solides pour lancer votre carrière en tant que Data Engineer.



le wagon

# Notre méthodologie

- ✓ 1 professeur(e) pour 7 étudiants
- ✓ Une classe à taille humaine
- ✓ Des professeurs passionnés pour vous aider
- ✓ Travaillez sur de vrais projets de startups
- ✓ Ajoutez votre projet final à votre portfolio
- ✓ Une pédagogie basée sur la pratique



"Si vous cherchez une formation qui offre un programme complet, un environnement d'apprentissage formidable et un incroyable soutien professionnel post-bootcamp, alors Le Wagon est fait pour vous !"



**Joseph Gulay**  
Désormais Data Analyst  
Ernst & Youngt



4.98 / 5  
2301 reviews



4.98 / 5  
2259 reviews



4.9 / 5  
714 reviews

# Choisissez le rythme qui vous convient le mieux

Sur campus ou en ligne, à temps plein ou à temps partiel !

## Temps plein

2 mois

Prêt(e) à plonger dans la Tech ? Rejoignez notre programme immersif. Du lundi au vendredi, de 9h à 18h.

- ✓ **Un diplôme en 2 mois**
- ✓ **Apprentissage en groupe toute la journée**
- ✓ **40 heures d'apprentissage par semaine**
- ✓ **Sur campus ou en ligne**

## Temps partiel

6 mois

Vous avez des contraintes d'emploi du temps ? Apprenez pendant votre temps libre grâce au temps partiel flexible.

- ✓ **Un diplôme en 6 mois**
- ✓ **Apprentissage en groupe, en ligne**
- ✓ **16 heures d'apprentissage par semaine**
- ✓ **En ligne ou sur le campus de Paris**

## Alternance

16 mois

Vous aimeriez vous former tout en étant rémunéré(e) ? Rejoignez nos formations 100% financées en alternance.

- ✓ **Un diplôme en 16 mois**
- ✓ **Apprentissage ou Professionnalisation**
- ✓ **3 mois de formation sur campus**
- ✓ **12 mois en entreprise pour pratiquer**



le wagon

# Le Wagon s'associe à Google pour former les futurs talents de la Data & l'IA

En rejoignant l'un de nos bootcamps Data, nos étudiants ont la possibilité d'intégrer le programme "Bootcamp Numérique IA" développé avec Google.

- ✓ Un accompagnement d'excellence
- ✓ Le consortium d'employeurs de Google
- ✓ Des certifications reconnues\*
- ✓ Votre formation 100% financée\*

[En savoir plus](#)



\*offre réservée aux étudiants éligibles

# Nos skill courses

40h

(nos formats courts)

Pour gagner en performance dans son job actuel ou découvrir un nouveau sujet.

✓ Temps partiel

✓ Session live en ligne

## Data Analytics

[Web Analytics & Tracking →](#)

[Data Analytics Essentials →](#)

## No-code / Low-code

[Growth & Data Automation →](#)

[Web Design & Webflow →](#)

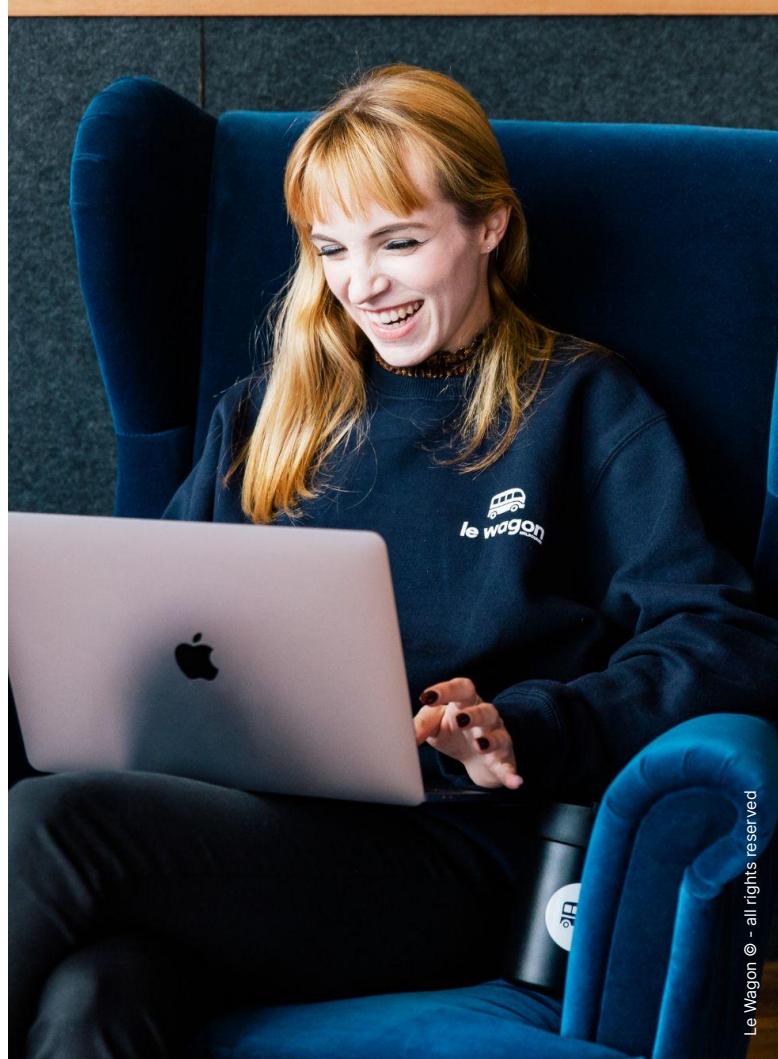
[Product Dev. & Airtable →](#)

## Data Science

[Python & Machine Learning →](#)



le wagon



# Tout le monde peut bénéficier de financement.

Maintenant que vous envisagez de faire vos premiers pas dans le monde de la Tech, vous vous demandez peut-être comment financer votre formation ?

## Options de financement

### ✓ Alternance

Formation **100% prise en charge** et **rémunération** pendant toute la durée du contrat (16 mois)



### ✓ Financements personnels

Paiement en **3 à 12 fois** avec Alma, prêt BNP Paribas



### ✓ Financements publics

Transition Pro, AIF, CPF, Aides régionales, POE!... **De nombreuses aides existent selon votre profil.**



## Tarif des formations

✓ Skill Courses : de 1490€ à 1690€

✓ Bootcamps : de 5900€ à 8900€

# Nos prochaines sessions

Sur le campus de Nantes

En ligne

## ✓ Développement Web (Temps plein)

- Du 15 janvier au 15 mars 2024 ❄️
- Du 15 avril au 14 juin 2024 🌸

## ✓ Data Analytics (Temps plein)

- Du 15 janvier au 15 mars 2024 ❄️
- Du 15 avril au 14 juin 2024 🌸

## ✓ Bootcamps à temps plein

- Du 29 janvier au 29 mars 2024 ❄️

*Développement Web, Data Analytics, Data Science*

## ✓ Bootcamps à temps partiel Flexible

- Du 10 février au 27 juillet 2024 ❄️

*Développement Web, Data Science*

## ✓ Formations courtes Skill Courses

- Du 23 octobre au 14 décembre 2023 🍁
- Du 15 janvier au 7 mars 2024 ❄️

# Ces entreprises recrutent nos alumni



le wagon

# Envie d'en savoir plus ?

- ✓ Votre situation / projet,
- ✓ Contenu des formations,
- ✓ Processus d'admission,
- ✓ Insertion professionnelle,
- ✓ Solutions de financement,
- ✓ Etc.
- ✓ Posez-moi toutes vos questions sur les formations du Wagon lors d'un RDV informel !
- ✓



le wagon



**Valentin Napoli**

Directeur de campus  
Le Wagon Nantes & Rennes

📞 06 79 99 24 96

✉️ [valentin.napoli@lewagon.org](mailto:valentin.napoli@lewagon.org)

linkedin [/valentinnapoli](https://www.linkedin.com/in/valentinnapoli)

# **Le Wagon vous offre un cours **gratuit** d'introduction à l'IA !**

Inscrivez-vous ce soir pour  
bénéficier de cette opportunité !



## **Introduction à l'intelligence artificielle**

- ✓ 8 heures de cours
- ✓ Outils gratuits uniquement



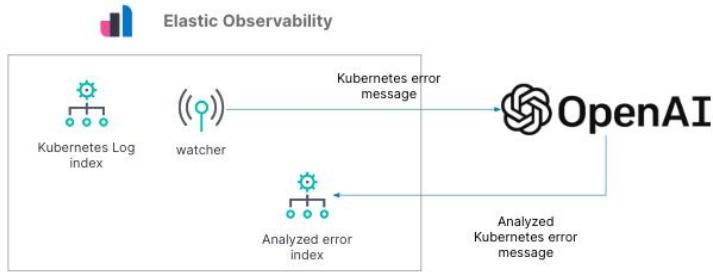
# Sponsors 2024

1y partnership

For  and 

Contact us ;)

# Observability: anomaly detection



👉 Elastic

Resources and Maybe Solutions ⓘ

AI SOLUTION EXPERIMENTAL

Hey there! 🌟

#### Problem Description

Looks like your app is hanging for at least 2000 ms! 😞 The issue happened at `-[NSBigMutableString getCharacters:range:]` in your app code. The app is stuck and unresponsive to user input, which is not good news! 💔

#### Proposed Solution

To fix this issue, you need to identify the root cause of the app hanging. You can start by profiling your app with a tool like Instruments to see if there are any performance bottlenecks. You can also review the code at `-[NSBigMutableString getCharacters:range:]` to see if there are any optimizations that can be made.

Once you've identified the root cause of the issue, you can make the necessary code changes to resolve it. The solution should be implemented in the affected method or class.

To prevent future occurrences of this issue, it's important to follow best practices for app performance optimization. This includes minimizing expensive operations, using background threads for long-running tasks, and avoiding blocking the main thread.

#### What Else

Hang in there! 🌟 App performance issues like this can be frustrating, but with a little bit of investigation and optimization, you can get your app running smoothly again!

Here's a hip hop rhyme to help you remember this error:

App hanging like a chandelier,  
2000 milliseconds of fear,  
Check your code and optimize,  
To make your app shine and rise!



Sentry

The screenshot shows the AWS CloudWatch Metrics interface for 'Log anomalies (6)'. It displays a list of detected anomalies with their priority levels (Low, High) and associated log patterns. To the right, there are four time-series charts showing the trend of anomalies over time.

Anomaly	Priority	Log pattern
47.6% increase in log event count	Low	[INFO] All good in loop <gt;> / <gt;> after waiting for <gt;> ms
3863.5% increase in log event count	High	REPORT RequestId: <gt;> Duration: <gt;> ms Billed Duration: <gt;> ms Memory Size: <gt;> MB Max Memory Used: <gt;> MB Init Duration: <gt;> ms
Unexpected pattern detected	High	[ERROR] Something happened in loop <gt;> / <gt;> after waiting for <gt;> ms
Unexpected pattern detected	High	START RequestId: <gt;> Version: \$LATEST
Unexpected pattern detected	High	REPORT RequestId: <gt;> Duration: <gt;> ms Billed Duration: <gt;> ms Memory Size: <gt;> MB Max Memory Used: <gt;> MB
Unexpected pattern detected	High	END RequestId: <gt;>

K8S 👈



K8SGPT  
KUBERNETES  
SUPERPOWERS

code size 273 kB build green release v0.3.3 openssf best practices passing Documentation

k8sgpt is a tool for scanning your Kubernetes clusters, diagnosing, and triaging issues in simple English.

It has SRE experience codified into its analyzers and helps to pull out the most relevant information to enrich it with AI.

Out of the box integration with OpenAI, Azure, Cohere, Amazon Bedrock and local models.

# Mixtral 8x7B & Mixtral 8x7B Instruct

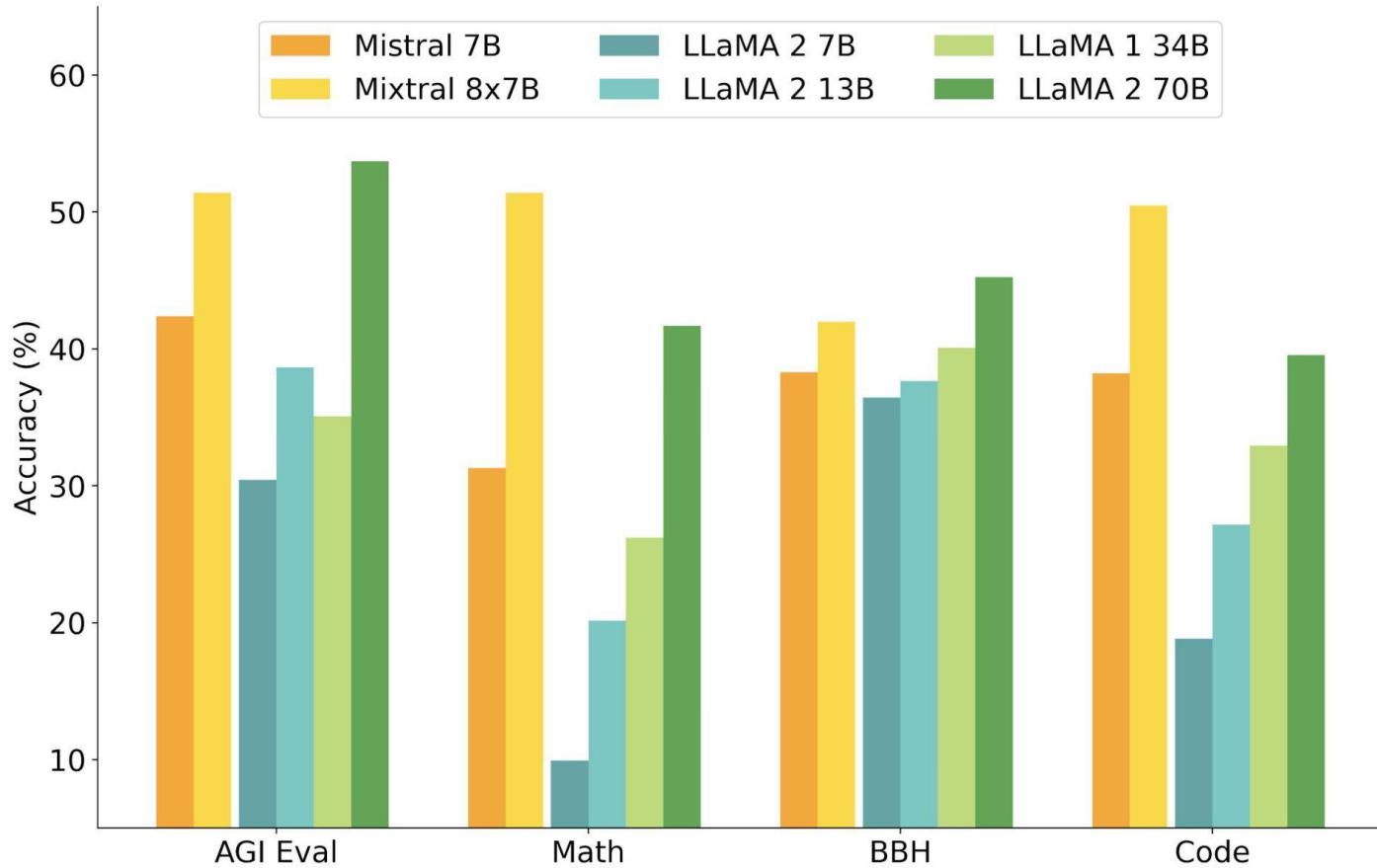
Mixture-of-Experts (MoE)

Outperforms Llama 2 70B and GPT3.5

32k tokens

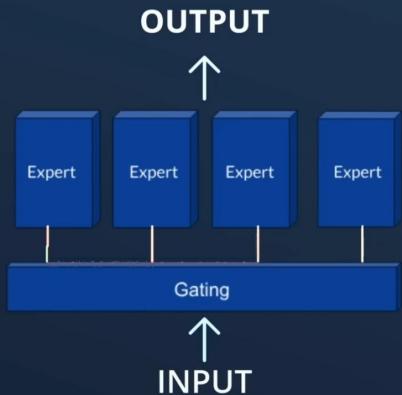
English, French, Italian, German and Spanish

	<b>LLaMA 2 70B</b>	<b>GPT - 3.5</b>	<b>Mixtral 8x7B</b>
<b>MMLU</b> (MCQ in 57 subjects)	69.9%	70.0%	<b>70.6%</b>
<b>HellaSwag</b> (10-shot)	87.1%	85.5%	86.7%
<b>ARC Challenge</b> (25-shot)	85.1%	85.2%	<b>85.8%</b>
<b>WinoGrande</b> (5-shot)	<b>83.2%</b>	81.6%	81.2%
<b>MBPP</b> (pass@1)	49.8%	52.2%	<b>60.7%</b>
<b>GSM-8K</b> (5-shot)	53.6%	57.1%	<b>58.4%</b>
<b>MT Bench</b> (for Instruct Models)	6.86	<b>8.32</b>	8.30



# Mixture-of-Experts (MoE)

## Mixture of Experts



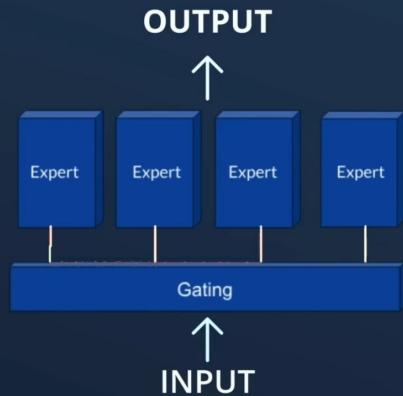
7B x 8 == 47B 😊

46.7B total parameters but only uses 12.9B parameters per token (2 experts).

Same speed and cost as a 12.9B model.

# Mixture-of-Experts (MoE)

## Mixture of Experts



More complex fine-tuning

Unequal load balancing

# Let's play a game: inference cost reduction



JJ — oss/acc @JosephJacks\_

...

Last week [@MistralAI](#) launched pricing for the Mixtral MoE: \$2.00~ / 1M tokens.

Hours later [@togethercompute](#) took the weights and dropped pricing by 70% to \$0.60 / 1M.

Days later [@abacusai](#) cut 50% deeper to \$0.30 / 1M.

Yesterday [@DeepInfra](#) went to \$0.27 / 1M.

Who's next ???

# Let's play a game: inference cost reduction



JJ — oss/acc ✅  
@JosephJacks\_

...

Last week [@MistralAI](#) launched pricing for the Mixtral MoE: \$2.00~ / 1M tokens.

Hours later [@togethercompute](#) took the weights and dropped pricing by 70% to \$0.60 / 1M.

Days later [@abacusai](#) cut 50% deeper to \$0.30 / 1M.

Yesterday [@DeepInfra](#) went to \$0.27 / 1M.

Who's next ??? 

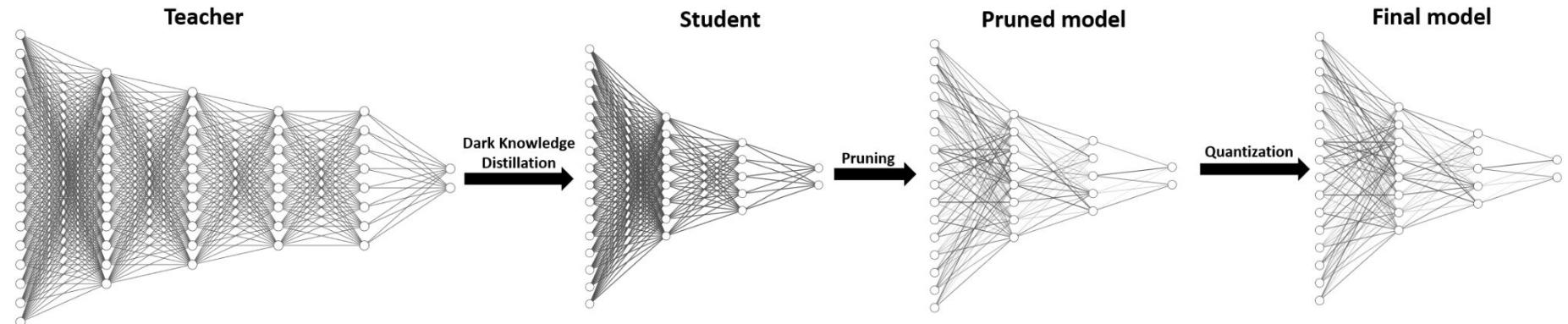
vs GPT-3.5-turbo cost: \$2/1M

# Mistral-medium on the road to GPT-4

## Chat Completions API

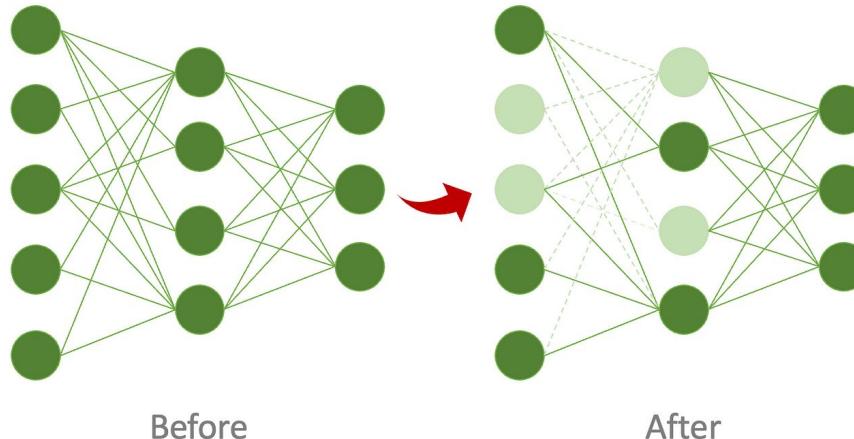
Model	Input	Output
mistral-tiny	0.14€ / 1M tokens	0.42€ / 1M tokens
mistral-small	0.6€ / 1M tokens	1.8€ / 1M tokens
mistral-medium	2.5€ / 1M tokens	7.5€ / 1M tokens

# Inference & compression

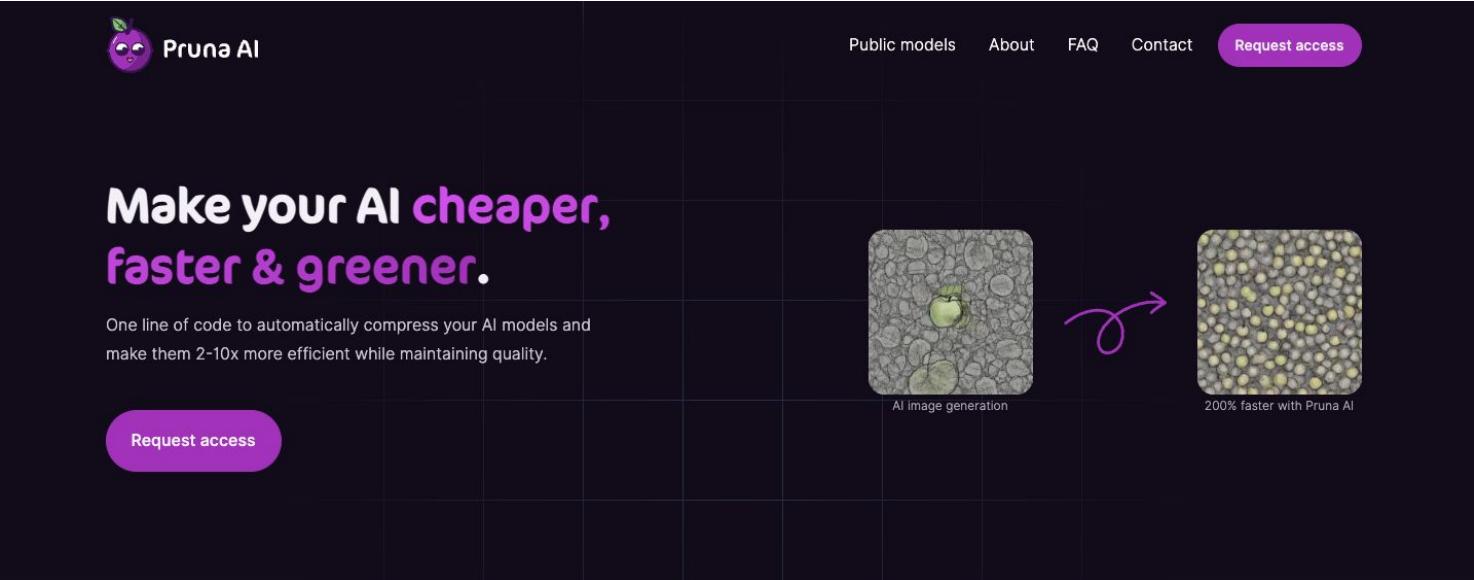


# Model pruning

Optimizing Deep Learning Models with Pruning  
A Practical Guide



# Model pruning



The image shows a screenshot of the Pruna AI website. At the top left is the Pruna AI logo, which features a purple apple with a neural network grid pattern. The top right contains navigation links: "Public models", "About", "FAQ", "Contact", and a purple button labeled "Request access". Below the header, a large call-to-action text reads "Make your AI cheaper, faster & greener." in white and purple. A subtext below it says "One line of code to automatically compress your AI models and make them 2-10x more efficient while maintaining quality." To the right of the text is a diagram illustrating model pruning. It shows two square grids of small circles. The first grid is labeled "AI image generation" and contains a single green circle. A purple arrow points from this grid to a second grid, which is labeled "200% faster with Pruna AI" and contains many green circles. At the bottom left of the main content area is another "Request access" button.

Pruna AI

Public models   About   FAQ   Contact   Request access

Make your AI cheaper,  
faster & greener.

One line of code to automatically compress your AI models and make them 2-10x more efficient while maintaining quality.

Request access

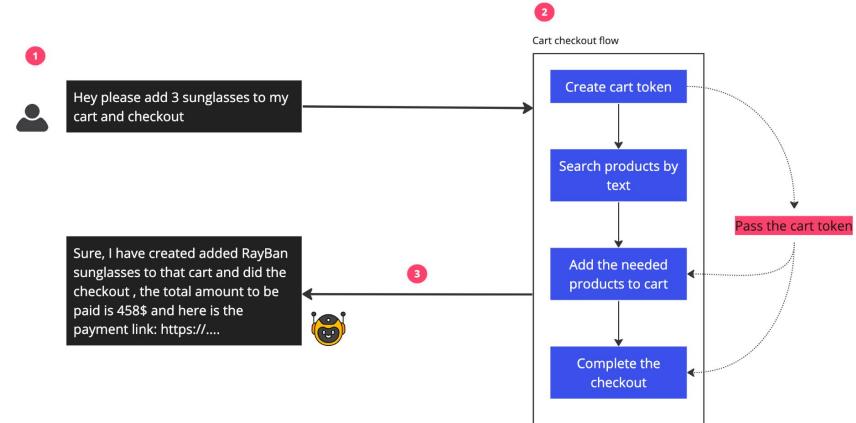
AI image generation

200% faster with Pruna AI

15% to 75% weight reduction - x2 faster

# OpenCopilot

<https://docs.opencopilot.so/introduction>



# Gladia: whisper-zero

Back



Jean-Louis Queguiner @JiliJeanlouis

...

Exciting news!

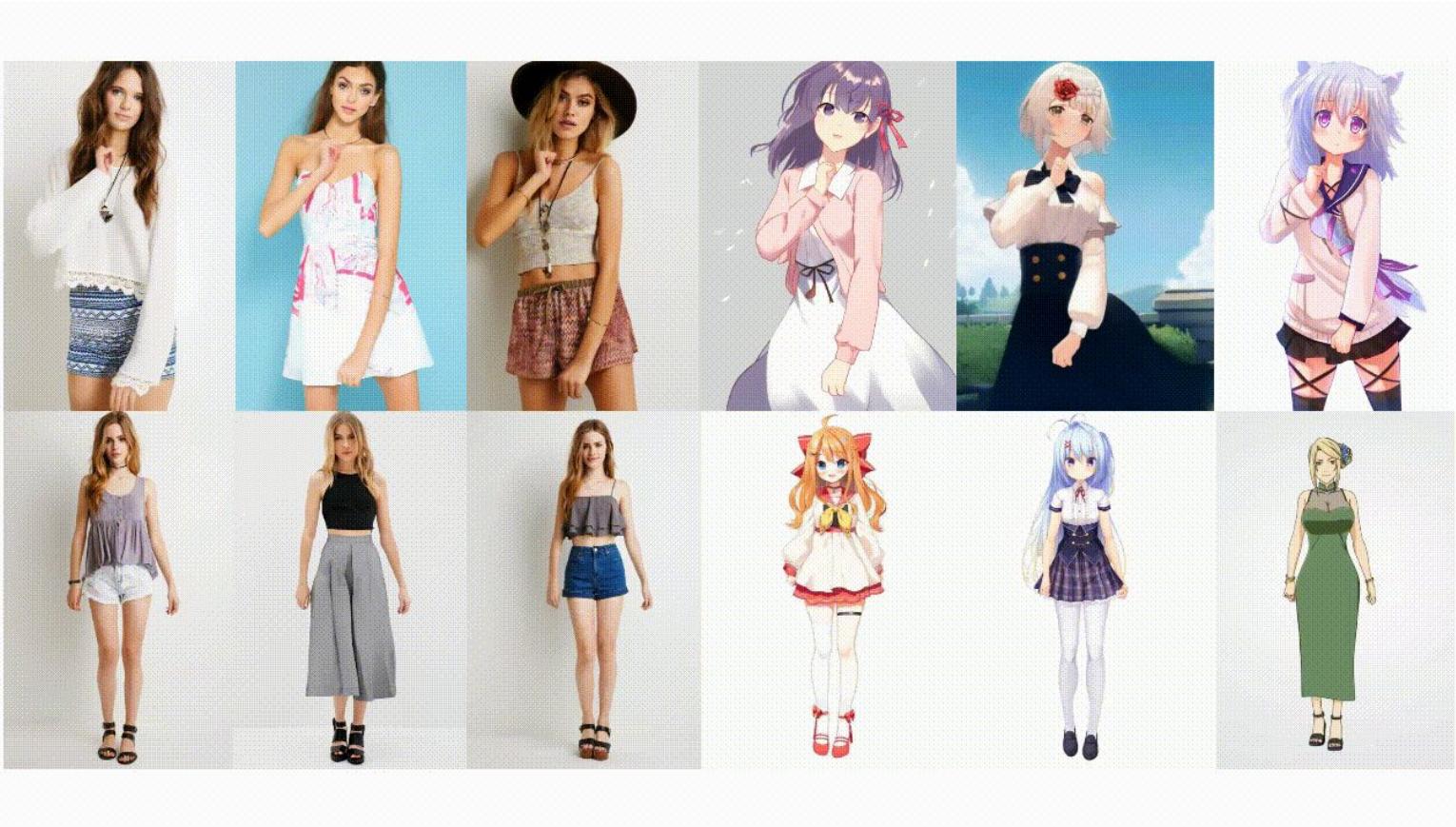
After months of hard work at [@gladia\\_io](#), we're thrilled to introduce Whisper-Zero, our latest model that solves major pain points.

We asked our customers what their biggest problems were:

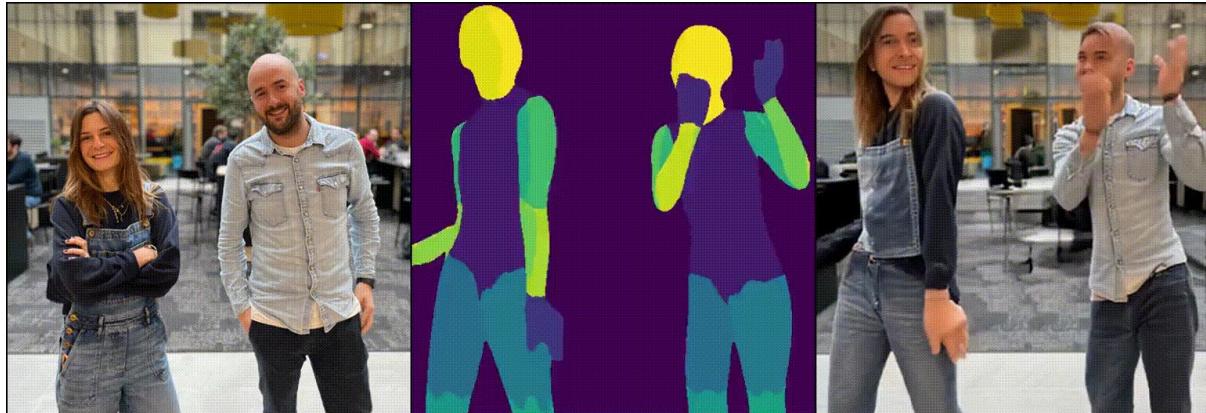
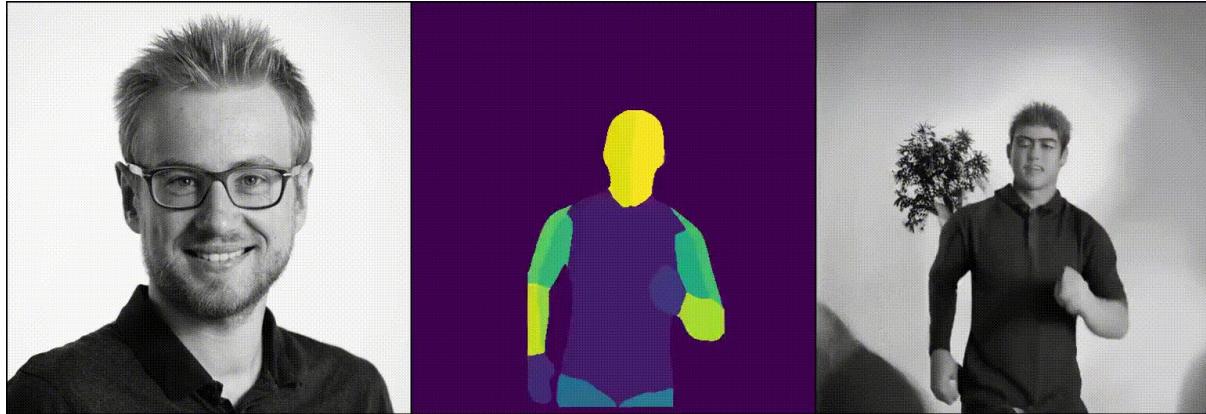
- **hallucinations** (which we fixed 0 hallucinations)
- **language detection** (which we highly improved in the 0's % error)
- **code-switching** (change the language on the fly 0 setup needed)
- **complex environment sensibility**
- **huge errors in numbers** leading to millions of \$ lost in our customers' critical business process automation

It leverages the Whisper architecture, extensively optimized on over **1.5 million hours of diverse audio recordings captured under various conditions, including hostile environments.**

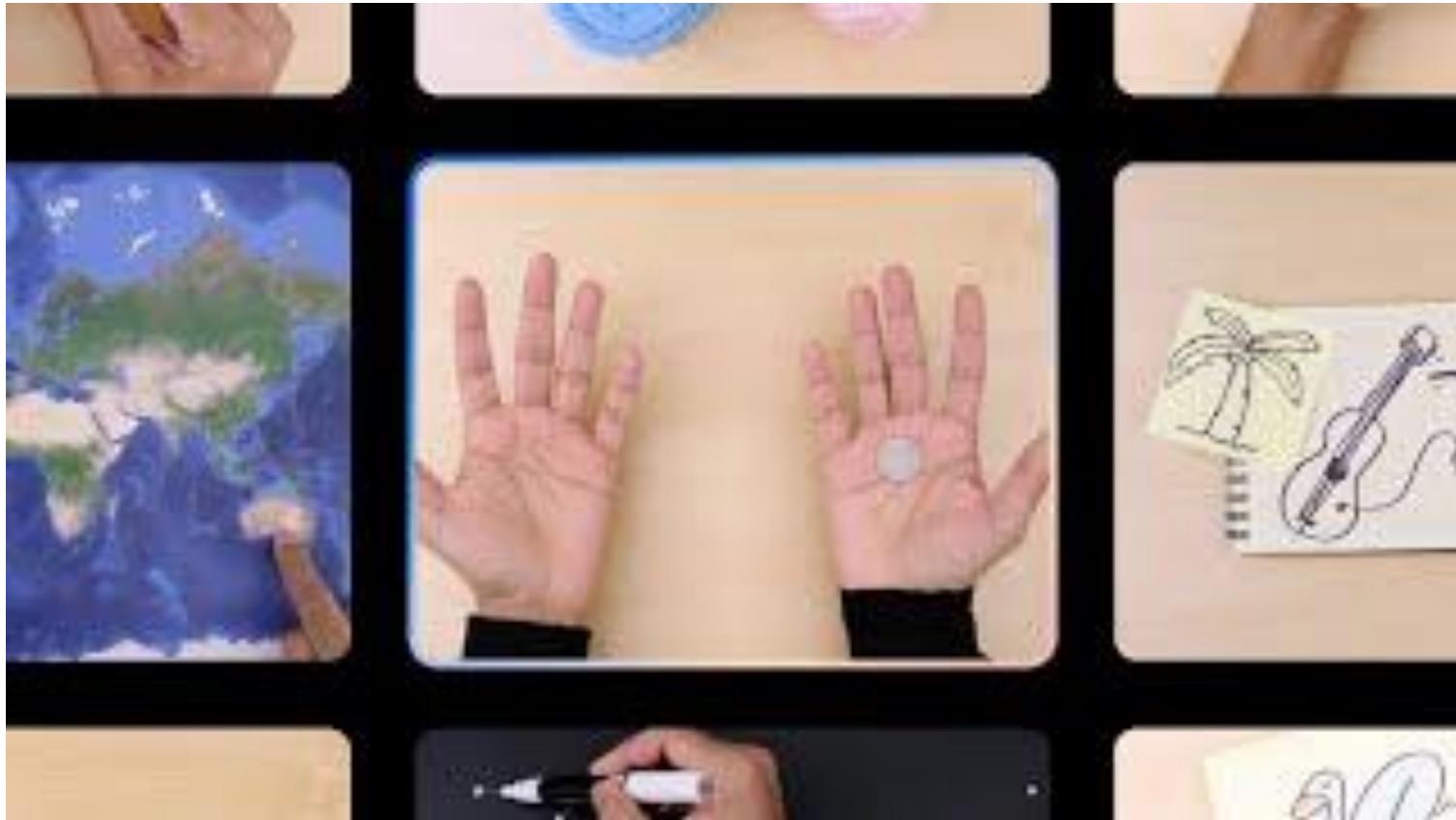
# About fake demo (Alibaba)



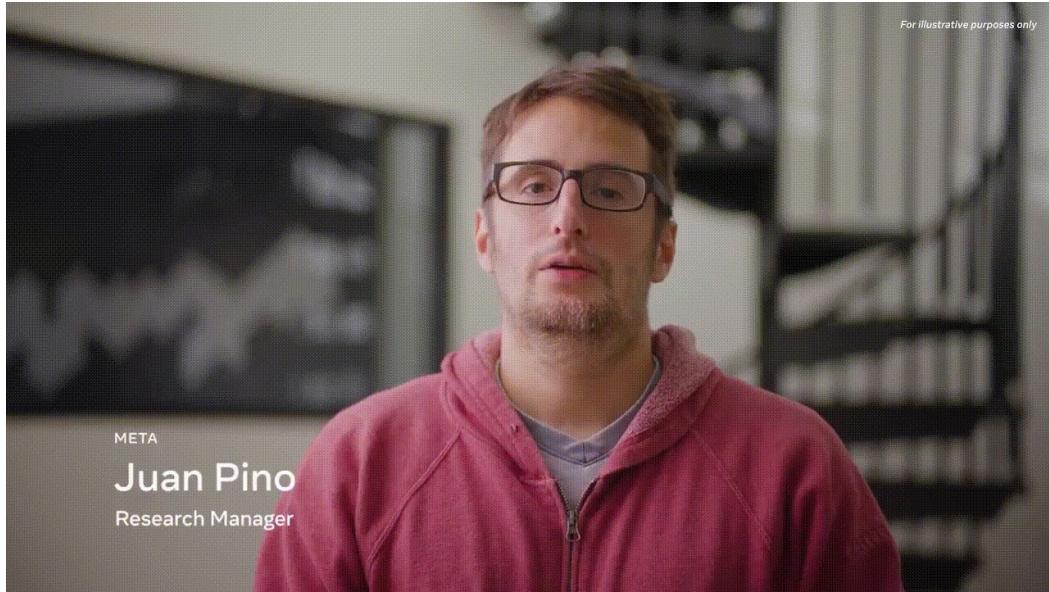
# About fake demo (Alibaba)



# About fake demo (Google Gemini)



# Real-time translation



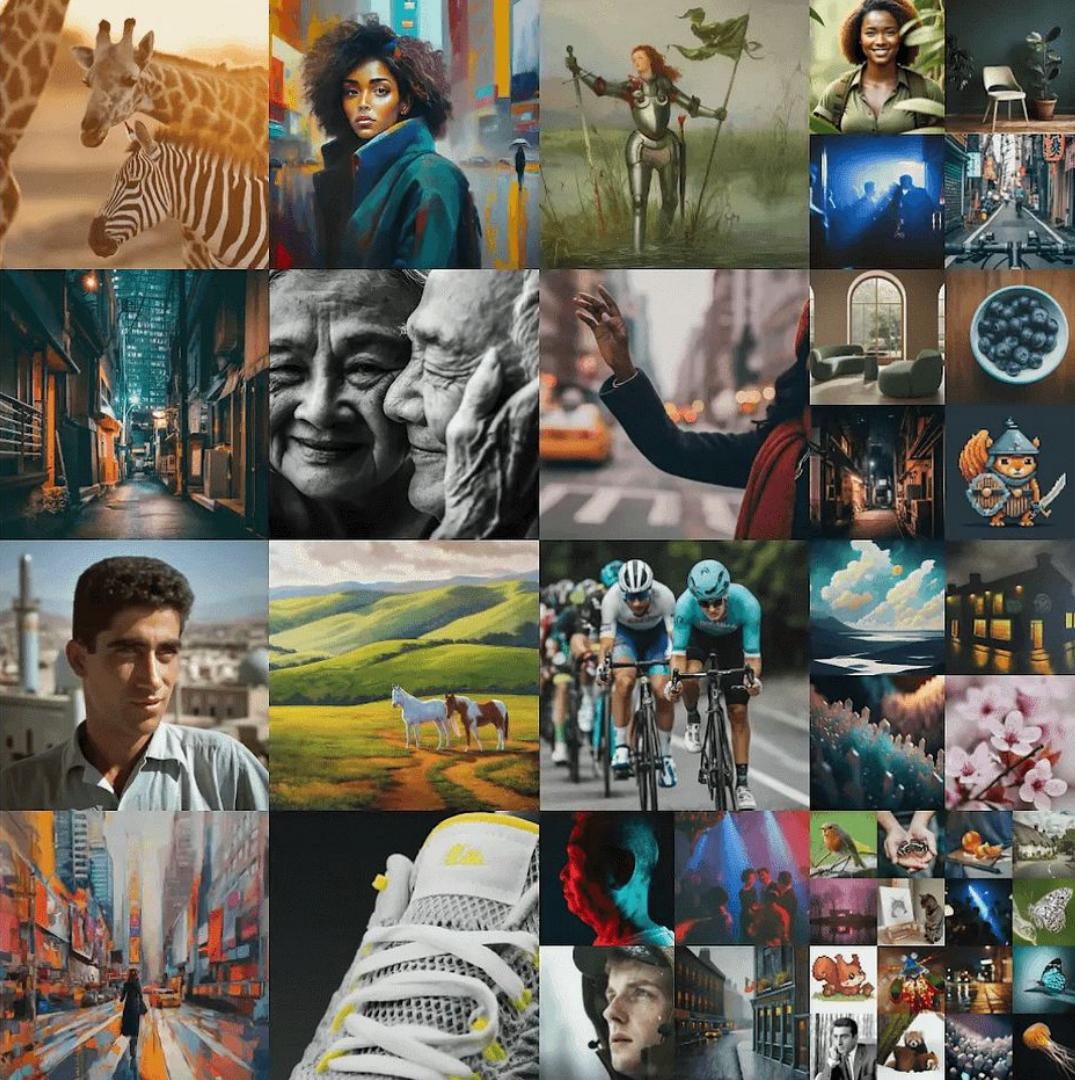
Real-time (2s latency)

Output before end of sentence.

Open-source weights

<https://x.com/BrianRoemmele/status/1735836569960173602?s=20>

# Imagen 2 most advanced text-to-image technology









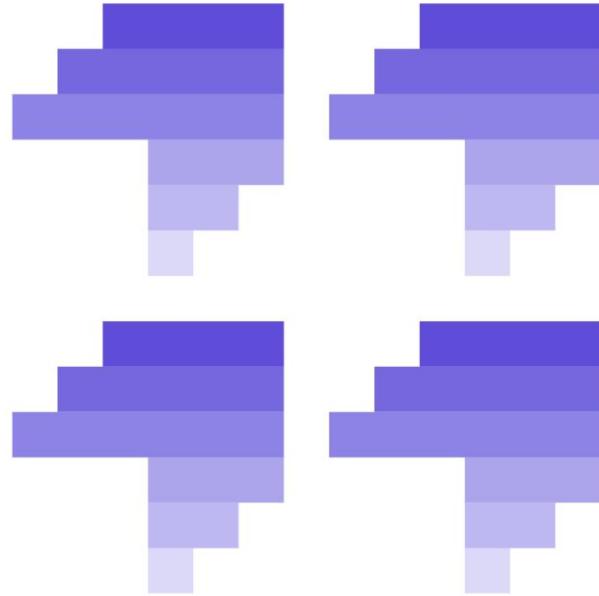


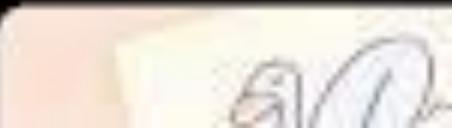
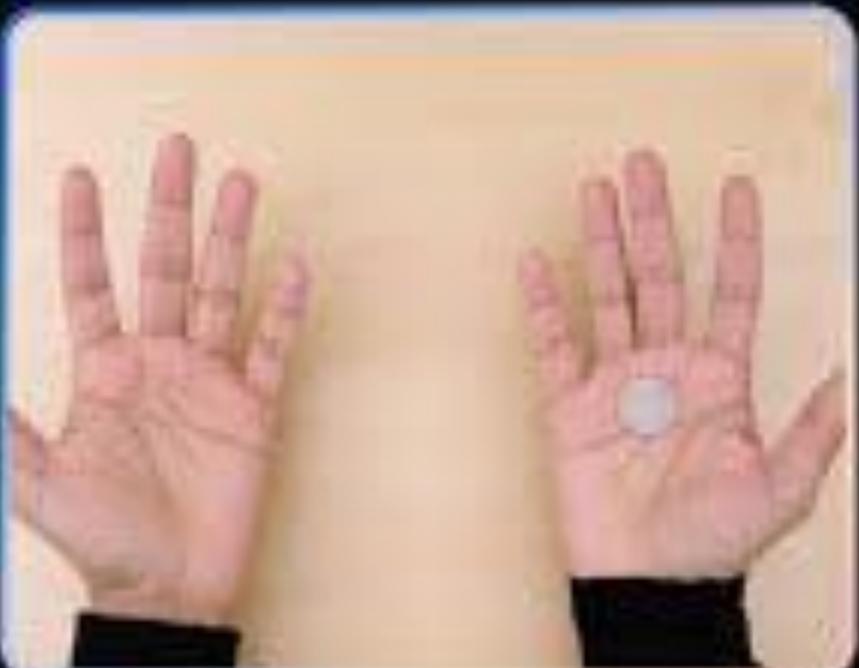
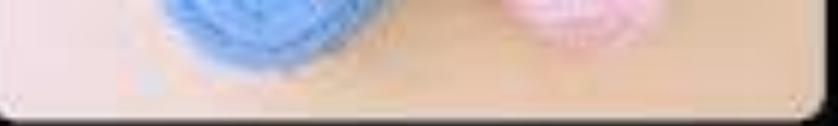
# MAGNIFIC AI



an island with palm trees



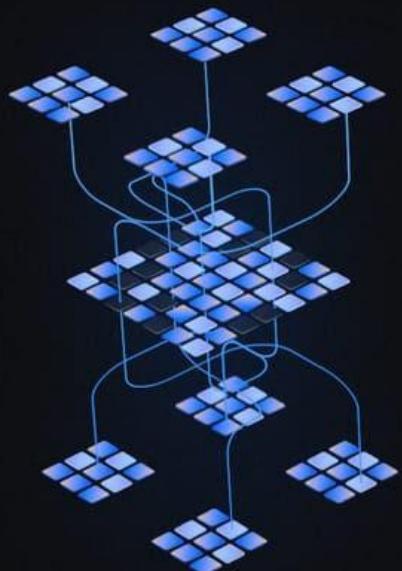




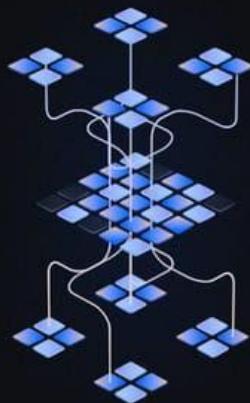
Stop



# Gemini



Ultra



Pro



Nano



GOAL GO  
THROUGH THE  
TALK TO  
AT

**Text:**  
+1 (314) 333-1111



**iOS App:**  
Pi, your personal AI



**Web:**  
[pi.ai](http://pi.ai)

The pi.ai logo, which is the word 'pi' in a bold, white, sans-serif font inside a light green rounded rectangular bubble.

**WhatsApp:**  
+1 (314) 333-1111



**Telegram:**  
[pi.ai/tg](https://pi.ai/tg)



**Instagram:**  
[@heypi.ai](https://@heypi.ai)



**Facebook  
Messenger**



# Pi



META

# Juan Pino

Research Manager

# Generate a TikTok Video from any content

1 Your video text

This is a text



Tip: use a link to a tweet, LinkedIn post, or blog post and we will magically retrieve the content  
⚠ If content is bigger than 500 characters, it will be summarized



2 Select Voice

Adam

american · deep · middle aged

0:00 0:05

Ryan Kirk - startup

English · middle aged · max

0:00 0:03

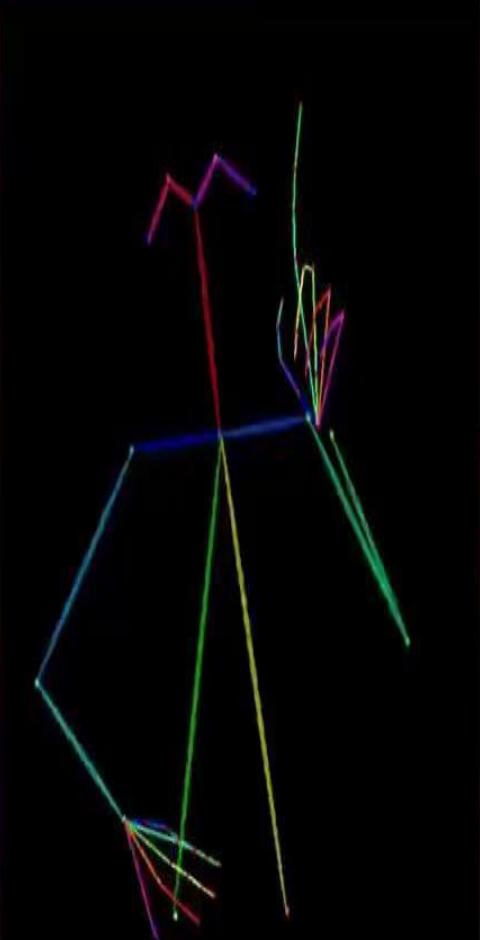
More voices ▾

Generate video

Start for free. No signup required.



Reference Image



Driving Pose



Ours





# MagicAnimate:

## Temporally Consistent Human Image Animation using Diffusion Model

Zhongcong Xu<sup>1</sup>, Jianfeng Zhang<sup>2</sup>, Jun Hao Liew<sup>2</sup>, Hanshu Yan<sup>2</sup>, Jia-Wei Liu<sup>1</sup>, Chenxu Zhang<sup>2</sup>, Jiashi Feng<sup>2</sup>, Mike Zheng Shou<sup>1</sup>  
<sup>1</sup>Show Lab, National University of Singapore   <sup>2</sup>Bytedance





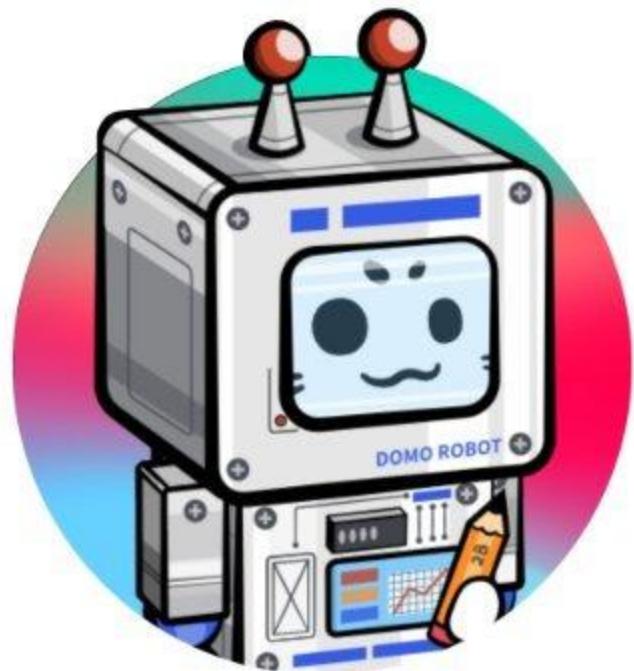
# Original





**RESEMBLE.AI**











Thierry Breton  
@ThierryBreton

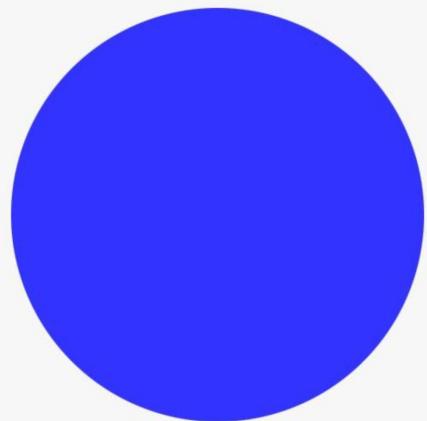


Deal!

#AIAct

### Continents that have an AI Regulation

- The EU
- Others



10:37 PM · 8 Dec 2023

1788 3762



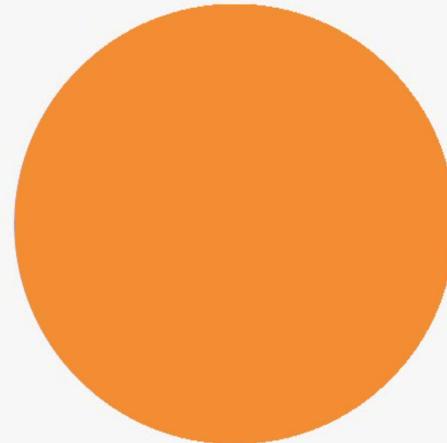
DIMENSION  
@DIMENSION\_YT



@ThierryBreton

### Continents that will benefit from the AI market

- The EU
- Others



2:30 PM · 9 Dec 2023

5 713

# #CONF'

*Barbara DELACROIX  
& Marvin SANT*

*Fondateurs @ Devana.ai*



*Le 19/12/23 à 19h*



*le wagon*

*"Le RAG : booster de connaissances pour la GenIA ?"*

