

Pruna AI

The AI Optimization Engine
cheaper, faster, smaller & greener AI

The AI race is moving too fast



Lost in AI evolutions



Lack of specialised talent



Too many AI models



Complex compute options

To be efficient & win your AI race you need partners with

PRODUCTION-READY TECHNOLOGY

Deep expertise in AI efficiency & reliability

BUSINESS ALIGNMENT

Flexible solutions compatible with evolving use cases

BUILDING TRUST

Accuracy and IP protection through sovereignty

We equip companies for their AI race

3rd-party AI APIs

Limited Open-Source Models

Limited Proprietary Model Development

Proprietary Model Development & Optimization

Self-Serve AI Enterprise Platform

Use our open-source models

AI Optimization engine

We're upgrading companies for the AI gold rush

ENTERPRISES THAT DEPLOYED AI IN 2023

42%

An additional 40% are stuck in experimentation. Top barrier is *limited AI skills* (32%) *High price* cited by 21%

AI SOFTWARE REVENUE 2025

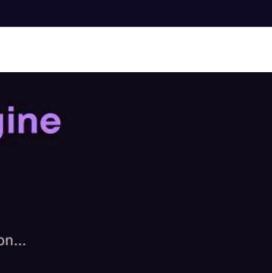
\$100B

GENERATIVE AI MARKET 2024

\$20B \uparrow 35% YEARLY

GENERATIVE AI MARKET 2030

\$110B



sources: AI Adoption Index, OMDIA, Grand View Research

Pruna is the AI Optimization Engine

2 lines of code for efficient inference

Combine AI efficiency methods
model pruning, quantization, hardware compilation...



Save time, money & carbon
for AI and Team productivity



Select your trade-offs
latency, memory, cost, energy...



Proven on 6500+ AI Models on Hugging Face (#1)

25%



LLM: GPU memory compared to original model

33%



LLM: Latency compared to original

50%



Speech to text: Latency compared to original

33%



LLM: Carbon emissions compared to original

Optimize Anything, Anywhere

NLP & LLMs
Llama 3, Mistral, Phi...

Image and Video Generation
Stable Diffusion, Consistency Models...

Computer Vision
ViT, MobileNet, Yolo, ResNet...

Audio
Whisper, Kaldi...



Compress
any hardware - any target metric

2 lines code

ML Engineers

AI-Based Agents & Assistants

Intelligent Document Processing

Process Automation

Customer Services Personalization

Predictive Maintenance

Fraud Detection

Clinical Research

Fit with Any Stack

DATA PLATFORM

STORAGE

PROCESSED DATA

Experimentation

DEVELOPMENT **TRAINING** **VALIDATION**

Production

OPTIMIZATION **SERVING**

German-French Leaders in Efficient & Reliable AI

Hiring thanks to our 6M seed a team of 15 people with ML Research Engineers.



BERTRAND CHARPENTIER
Cofounder & Chief Scientist

PhD in ML at Technical University Munich
ML MSc KTH & Ensimag, summa cum laude
ML Research at Twitter, Stanford & Telecom

Leading expert in ML efficiency & uncertainty
Created scikit-network with >500k downloads

LinkedIn - Google Scholar



STEPHAN GÜNNEMANN
Cofounder & Chief Strategy Officer

Professor of ML of Technical University Munich
Head of ML research group with > 30 PhDs
Partnerships with Google, Siemens, BMW...

Leading expert on reliable machine learning
Work productized on > 100 million users

LinkedIn - Google Scholar

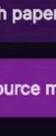


RAYAN NAIT MAZI
Cofounder & CEO

MEng Centrale Paris, MPhil HKUST
3rd time cofounder
AI product builder

Leading expert on low-code product dev
RAG research copilot, experiment manager tool
for biotech, EdTech chatbots in Indian slums

LinkedIn



JOHN RACHWAN
Cofounder & CTO

ML MSc at TUM, summa cum laude
ML Engineer in Efficient ML at TUM
AI Team Lead at Design AI (acq. by Helsing)

Won EDA Defence Innovation prize 2021 by
leading his team in the full reimplement of
DeepMind's AlphaStar in military simulators

LinkedIn

We Are Trusted and Mature

>270+ research papers

>6.7k open-source models

>200k monthly downloads

Backed & Supported



Pruna AI