# PALI

## Performant AI
## Launcher for Inference

**Download the model with a single click**

**One more click to launch your own service**

## Struggling with the setup to run your AI models?

Public institutions / Internal service operators
**"Need easy access to models within organization users."**

Startups / Service planners
**"Generative AI APIs are too pricey and complicated."**
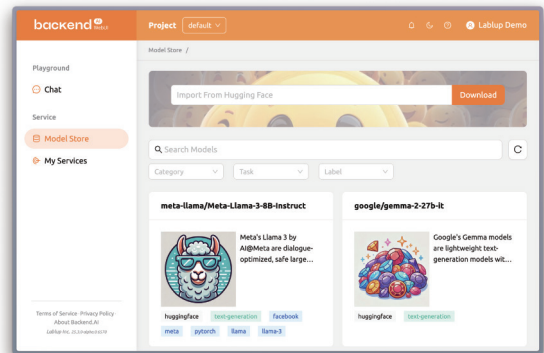
Developers / Researchers / AI beginners
**"GPU Setup? Docker? Configuration is way too complex."**

## Start your AI journey with PALI in just one click.
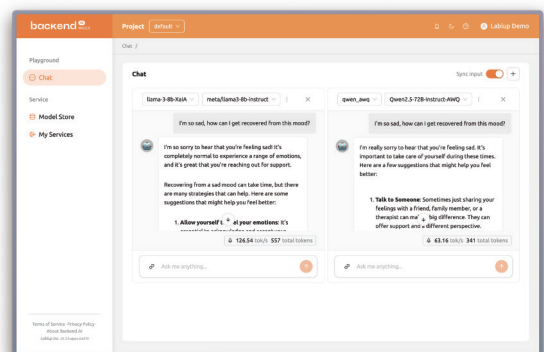
lablup

# Model Store
## Various models in one place

- Provides curated catalog of public GenAI models
- Supports importing models from Hugging Face
- Supports NVIDIA NIM™ Inference Microservices
- Allows uploading user custom models
- Chat UI for instant model testing
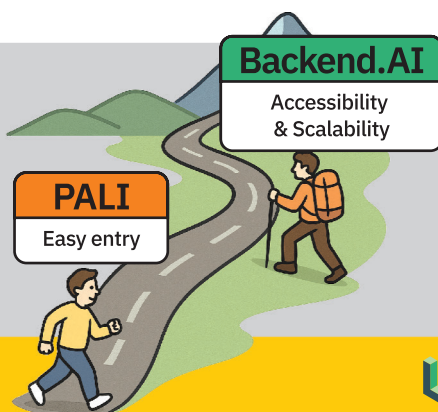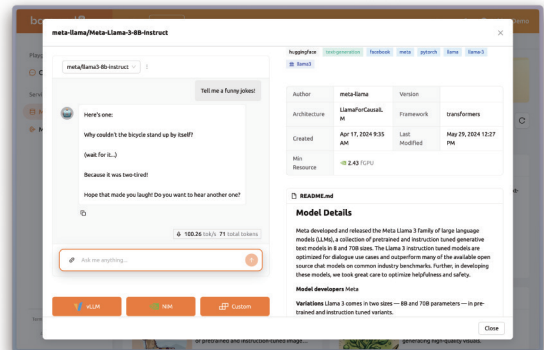
# Playground
## Instantly run selected models

- Test language models tailored to your service
- Fuse: combine multiple outputs for the best answer
- Compare models side by side with Multi-chat UI
- Supports image-based multimodal input/output
- Create, save, and reuse custom prompts

# Model Player
## Instantly run selected models

- Click & Play: Fast model execution
- Supports inference runtimes like vLLM / NIM / Hugging Face TGI
- Provides endpoints to integrate models with external services

**Backend.AI**
Accessibility & Scalability

**PALI**
Easy entry

PALI is your starting point for GenAI — launch complex AI setups in just two clicks. Works seamlessly across on both on-prem and cloud environments.

New to AI? Start simple with PALI.
Power user? Go deeper with Backend.AI.

**Start your AI journey with Lablup.**

lablup

contact@lablup.com
https://www.backend.ai
https://github.com/Lablup/backend.ai