

LGTM.  Base Architecture by  
RWKV

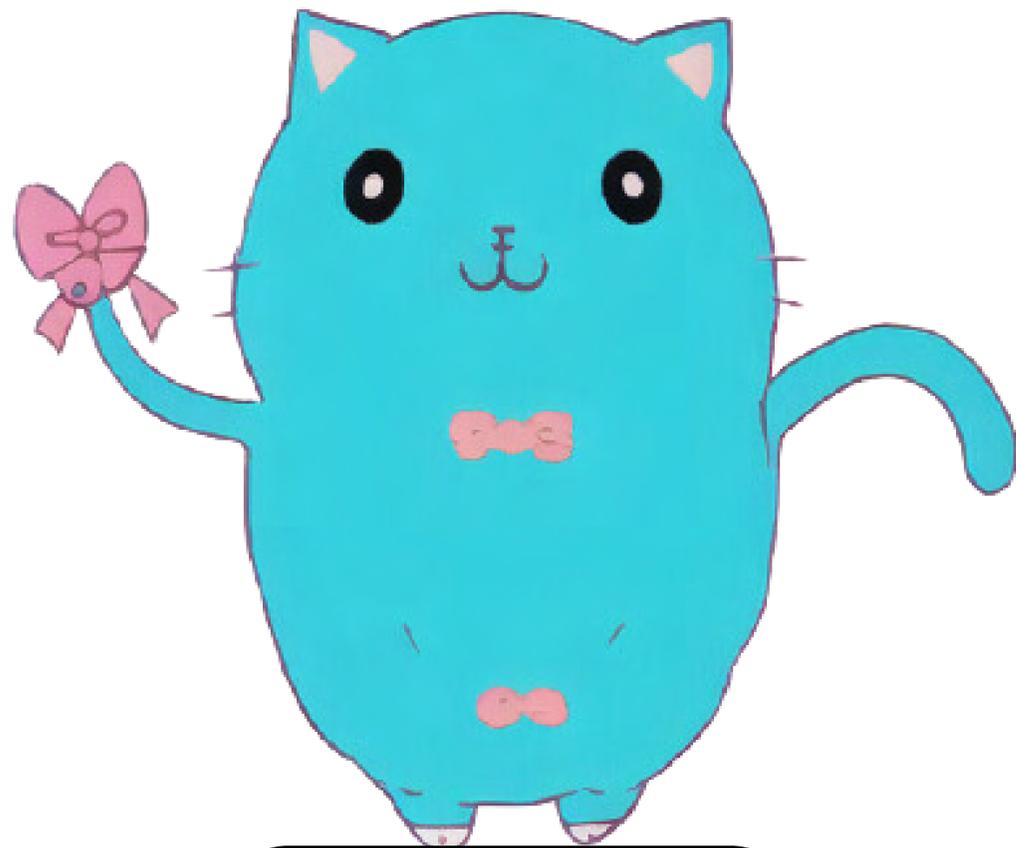
Language Gateway to Tomorrow Module

みんなが使えるAIを。

**LGTM.**  Base Architecture by  
RWKV  
Language Gateway to Tomorrow Module

# About LGTM

# LGTM開発統括よりご挨拶



OpenMOSE



**OpenMOSE** @\_m0se\_ • XX月XX日

LLMが数多く誕生していく中でLGTMはどこまでも小さくて賢い言語モデルを追求していきます。

どんなスマートフォン / PCでもAIが当たり前機能の一つとして備わっているので、小さくて賢いモデルが必要になると考えております。

「身近なAI」を築き上げる日本の先駆者として日々LGTMを磨き続けます。

詳しくはX(@\_m0se\_)をご覧ください。

❤️ 10100

# LGTMって何？

大規模言語モデルの一つ、RWKVを基にした言語モデルです。

「みんなが使えるAI」を実現するために、誰もが触りやすく、製品を作りやすい言語モデル LGTM を開発しました。



**誰でも安く簡単に**  
**オリジナルAIサービスの開発が可能に**

# お悩みの声

社内規則を社員に周知したい、、  
社員同士の会話を増やしたい！

誰にも話せない愚痴を聞いてくれる  
相手が欲しい

自分の代わりにコーディングしてくれる人  
がいたら、、

AIを取り入れたいけどクラウドだと  
情報漏洩が怖い、、

**ChatGPTやLLaMAを使っても出来ない、  
AIの活用方法が分からない、というお悩みの声をいただきます**

# LGTMの活用シーン



## 社内の情報屋 wikiボット

社内のあらゆる情報を学習させることで社内ルールや、社員の趣味や嬉しいことなどをいつでも聞いて相談ができ、コミュニケーションの促進だけでなく新卒社員の退職率軽減にも役立ちます。



## コーディングAI

プログラマー、足りてますか？御社のプログラミング規則、独自ライブラリ、これまでに制作したコードを学習させることで、御社独自のコーディングAIを作成することが可能です。



## 自分だけの理想の話し相手

自由に口調や性格をカスタマイズ可能なキャラクター性の高いAIを作ることができます。AI Tuberの制作や、介護現場での話し相手、また、あなただけの理想のパートナーを作ることが可能です。

# LGTMは他にも様々な課題を解決することができます

# 比較表

|            | LGTM  | GPT   | LLaMA   |
|------------|---|---|---|
| 導入レベル      | 初心者～上級者   | 中級～   | 上級者   |
| 生成<br>スピード | <br>(RTX3090にて7B 200Token/sec) |  |  |
| 学習         | クラウド/オンプレミスで簡単<br>最短2分で学習   | 限定的なファインチューン/<br>基本はプロンプトチューニング   | 難しい   |
| コスト        | 世界最安・月額固定でファイン<br>チューニング可能  | APIによる従量課金  | 高額な機材が必要  |
| ビジネス面      | 製品用特化モデルの作成にかかる<br>時間が削減できるため、エン<br>ド製品を早くリリースできる   | OpenAI社に依存したシステムを<br>構築する必要がある  | ファインチューニングに時間が<br>かかるため、製品実用化までに<br>時間がかかる  |

# 直感的に操作できる学習環境 LGTM Easy Trainer

「LGTM Easy Trainer」では、所定の形式のCSVファイルを読み込ませて、画面の指示に従っていただくだけで、あなただけの特化言語モデルを作成できます。

直感的な操作

簡単なデータセット

圧倒的な学習スピード

The screenshot displays the LGTM Easy Trainer interface. On the left is a sidebar with navigation options: 'LGTM for Business Console', '学習する / Easy Trainer', 'テストする / Test Model', and '運用する / Inference Server'. The main area is titled 'LGTM Easy Trainer' and contains three numbered steps:

- 1. 学習に使用するデータセットをアップロードしてください**  
Upload your dataset csv file  
ここにCSVファイルをドラッグするか、クリックしてファイルを選択  
Drag and drop file here (Limit 200MB per file • CSV) [Browse files]  
CyberloungeQA\_dataset.csv 23.7KB
- 2. 学習に使用するVRAMサイズを選択してください**  
Select the amount of VRAM usage for training  
 12GB  16GB  24GB
- 3. 学習するモデルのパラメータ数を選択して下さい**  
Select the number of parameters of the model to be trained  
 0.4B  3B  7B

A red button at the bottom of the steps says 'トレーニングを開始 / Start Training'. On the right side, there is a 'Loss' graph showing a fluctuating line over 200 steps. Below the graph, the 'Epoch: 0' progress bar shows 'step: 15/200 Time remaining: 07:08'. A 'Terminal' window displays the following log output: 'Epoch 0: 7% | 15/200 [01:00:22<07:08:22, 2.25s/it, loss=1.045012944768353, lr=2.45e-6, REAL it/s=0.547, Kt/s=2.240]'. There are dropdown menus for 'Log' and 'Log(RAW)'.

# LGTM 技術者向け説明

## 小規模で 高い日本語性能

LGTMでは、30億パラメータ  
および70億パラメータのモ  
デルを提供予定

少ない計算資源で高い効率  
と性能を追求した日本語の  
理解と生成に特化させた  
モデル

## コンテキスト数 無限処理

多数の言語モデルが採用し  
ているトランスフォーマー  
アーキテクチャとは異な  
り、LGTMはRNNを使用

時系列データを扱うのに優  
れており、自然言語の流れ  
をより連続的に捉えること  
が可能

## RNN型 アーキテクチャ

RNNの特性として、入力の  
長さに依存しないため、非  
常に長い文脈も理論上処理  
可能

長い文書や会話の理解に有  
利であり、コンテキストが  
非常に重要な翻訳や要約タ  
スクにおいて強みを発揮

## 高速な推論

RNNアーキテクチャのメリ  
ットとして、RNNは過去の  
情報を隠れ状態として保持

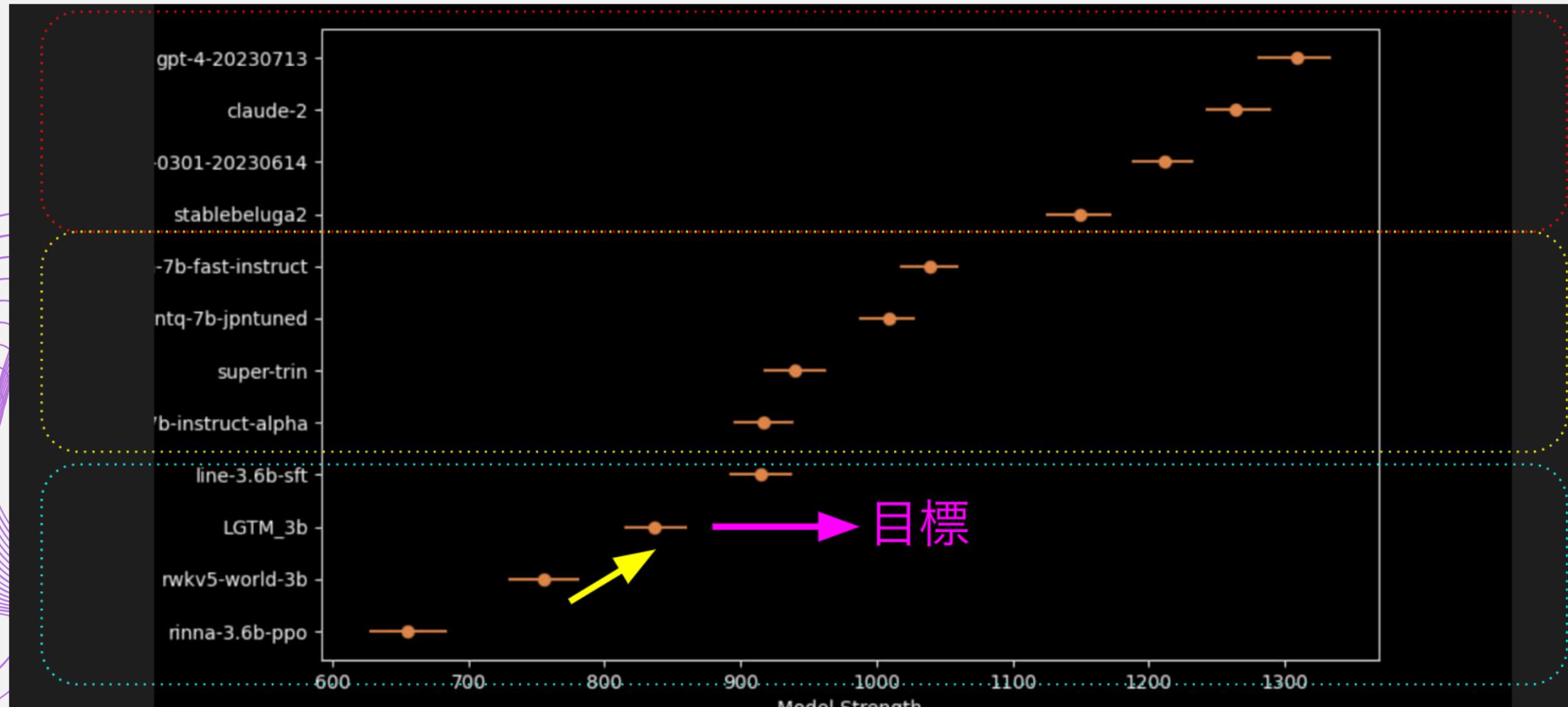
Transformer系と比べ再計  
算頻度を低減し、少ない  
GPUもしくは、CPUでの高速  
な推論が可能

# 日本語性能検証結果

## LGTM 3bの日本語性能(RAKUDAベンチマーク)

検証日：2024年3月19日時点

Transformer系と遜色がないレベルに到達



大規模  
70b+領域

中規模  
7b領域

小規模  
3b領域

# 製品化までの流れ



## ファインチューニング

- 1 少量のターゲットデータで高度に最適化されたモデル調整
- 2 ユーザー固有の用途や業界に特化
- 3 短時間でのデプロイメントが現実



## RAG

- 1 関連情報の検索と内容生成が一体で行われる
- 2 データベースからの情報取得と自然で正確な文の生成



## 製品例:チャットボット

- 1 自然な日本語で流暢に対話し、ユーザーからの質問に対して適切な回答
- 2 幅広いトピックに対応した満足度の高い対話体験

# 技術資料 (LGTM for Business)

## システム稼働要件(推奨)

OS: Ubuntu 22.04 LTS  
 CPU: x86-64互換CPU  
 RAM: 128GB以上  
 GPU: 90TFlops(BF16), VRAM 24GB以上  
 CUDA 12.3以上 もしくは ROCm 6.0以上  
 が動作するもの

## 配布形態

実行環境はDockerイメージで配布  
 (導入後は自動更新されます)

※ インターネットに接続できない環境への導入につきましてはお問い合わせください。

## サーバ要件

- SSHでアクセスできる状態であること
- Webサーバーが動作し、クライアントとの送受信ができる環境であること

## 推論結果 (テキスト生成) の利用方法

LGTM for Businessを導入したサーバから発行されるAPIエンドポイントとAPIキーをクライアント側で設定してください。  
 (詳細はオンラインダッシュボードからご確認ください。)

## 学習方法

オンラインダッシュボードの「LGTM Easy Trainer」から学習 (ファインチューニング) が可能です。

# 技術資料 (LGTMのモデル構造・特徴)

## モデル構造

LGTMはRNN(回帰型ニューラルネットワーク)を使用したLLM(大規模言語モデル)です。

基本モデルアーキテクチャとして、RWKVという言語モデルを採用し、それを元に自社開発の学習手法とデータセットを用いてチューニングを行っています。

## モデルの特徴

- ・ 実行およびトレーニング時の計算が少なく、大きなコンテキストサイズを持つTransformer系モデルと比べて、計算能力の要求が10分の1以下。
- ・ 計算負荷が任意のコンテキスト長に線形でスケールする。(Transformerは二乗にスケールする)
- ・ ベースモデルが多言語でトレーニングされており、日本語だけでなく英語、中国語にも対応可。

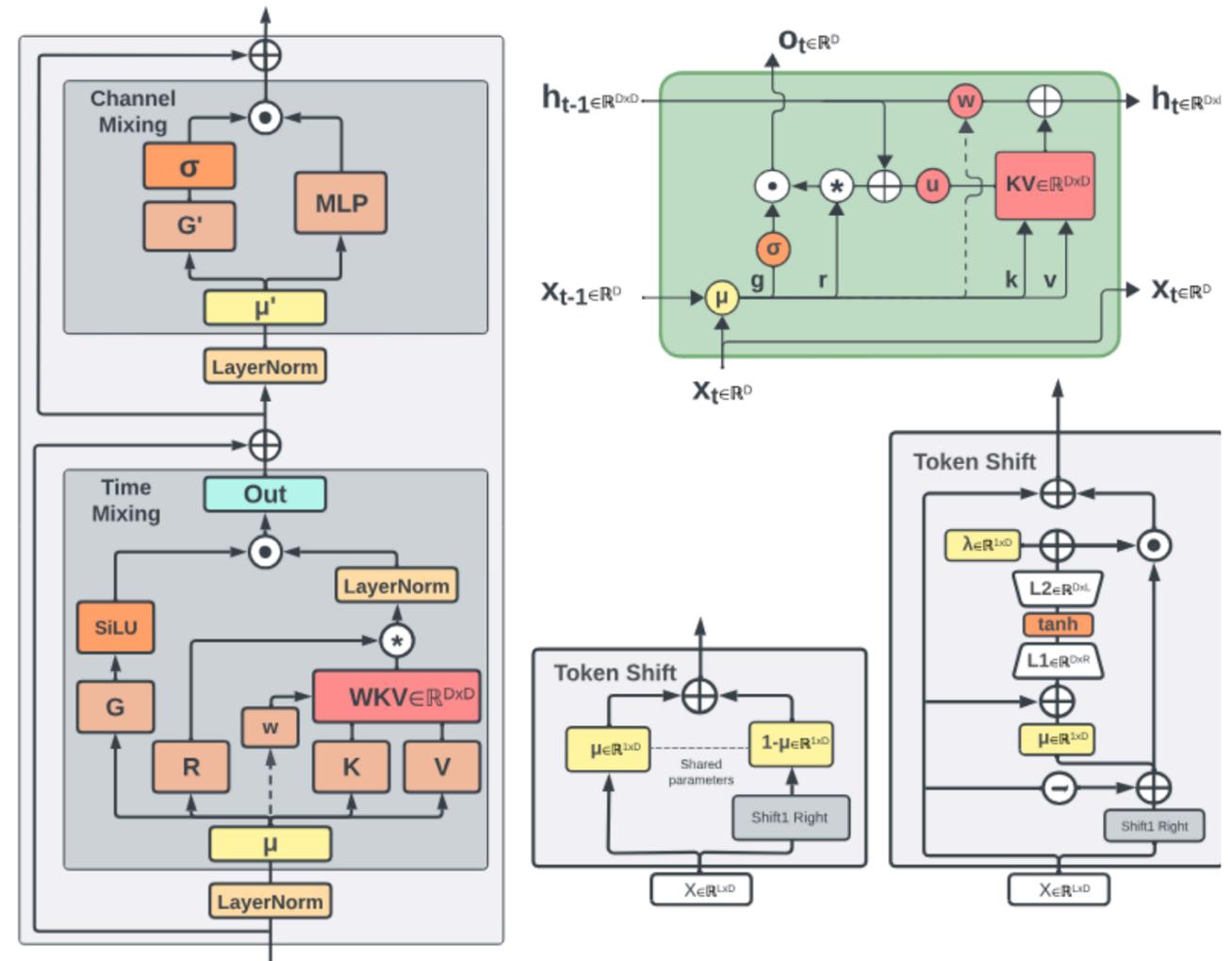


図. モデル構造

Eagle and Finch: RWKV with Matrix-Valued States and Dynamic Recurrence (<https://arxiv.org/abs/2404.05892>)

# 技術資料 (学習手法)

## 学習手法(Fine-Tuning)

「LGTM Easy Trainer」では State-Tuning という手法と RLHF (人間からのフィードバックによる強化学習) を組み合わせたものを採用しています。

その他の手法を使用してLGTMをファインチューニングする方法につきましてはお問い合わせください。

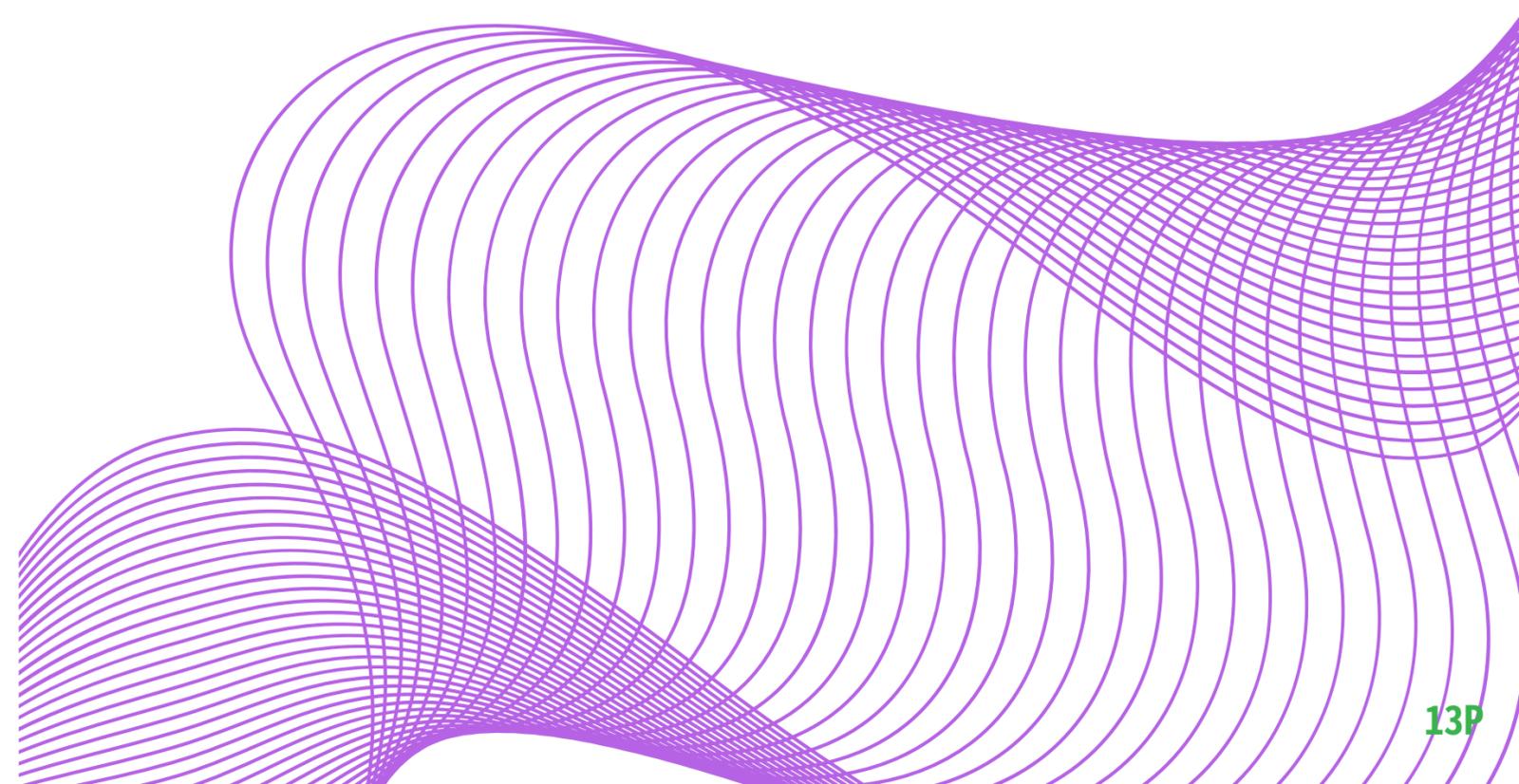
## State-Tuning

State-Tuningとは、入力データセットに対して適切な promptベクトルをバックプロパゲーションし、その時点のモデルのニューロン発火状態をニューラルネットワーク内の各層に埋め込む手法です。

## RLHF

RLHFとは、機械学習モデル学習時に人間の価値基準を元に、趣向に沿う回答が出るように調整する強化学習の手法です。

本製品では弊社独自実装の学習手法によりこれを実現しています。



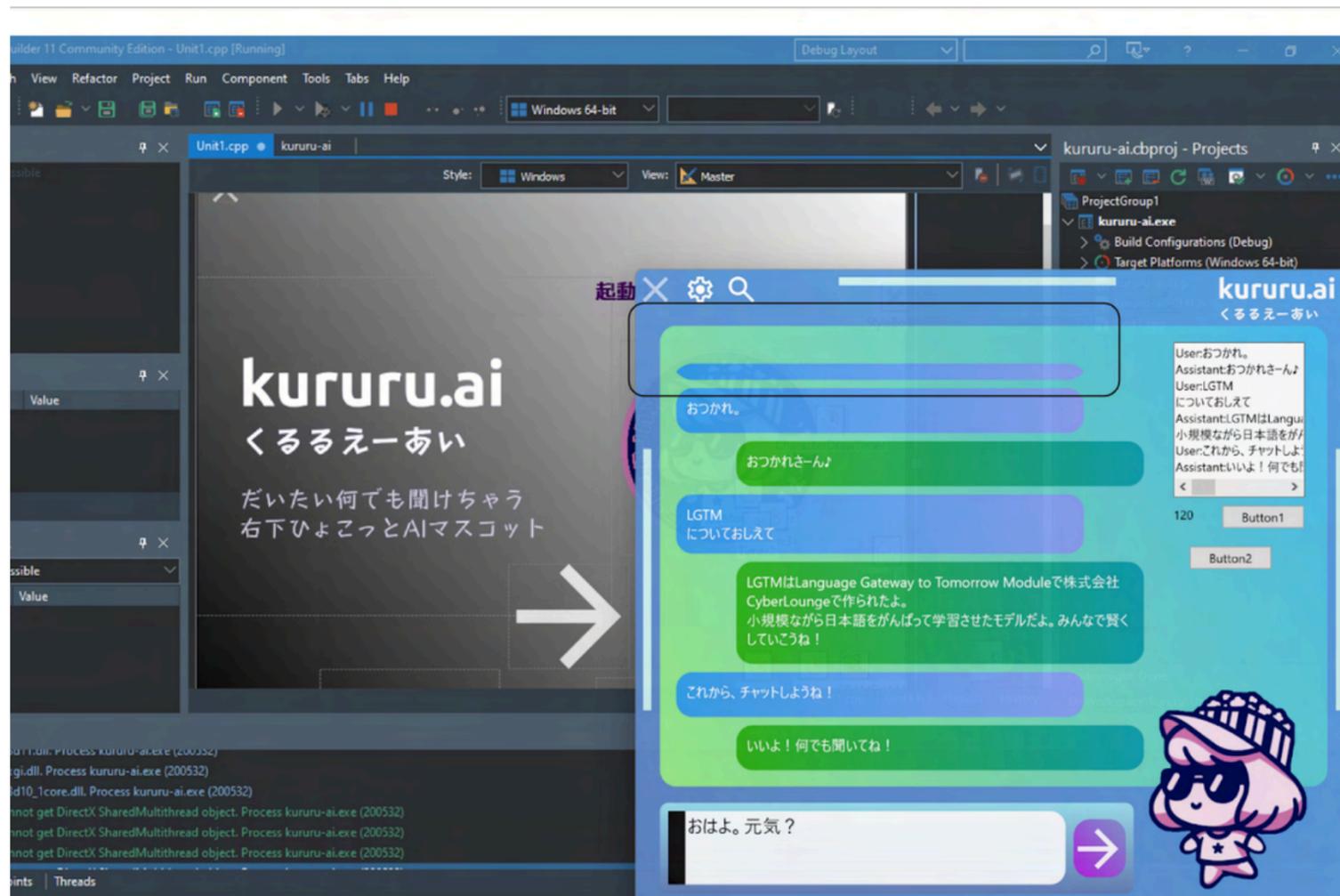


# kururu.ai

AIデスクトップマスコット

# kururu.aiは、 LGTMを使った無料AIチャットボットです。

ネット環境不要の  
Windowsアプリ



1 PCにインストールするだけ  
オフラインで利用可能

2 音声合成や天気予報など  
プラグインの追加が可能

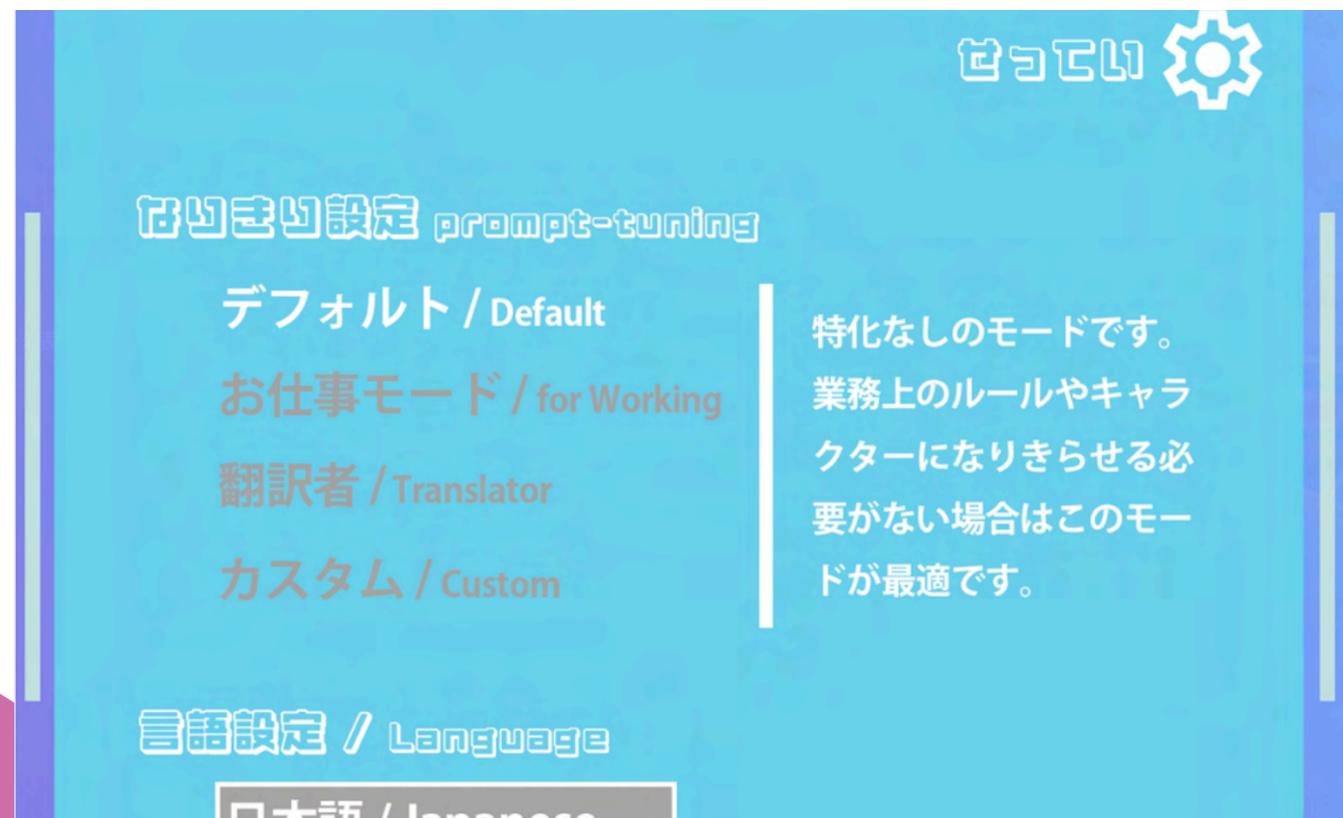
3 プロンプトチューニング  
でキャラ設定ができる

4 クオリティを上げるなら  
ファインチューンで簡単学習

# 簡単にキャラ設定

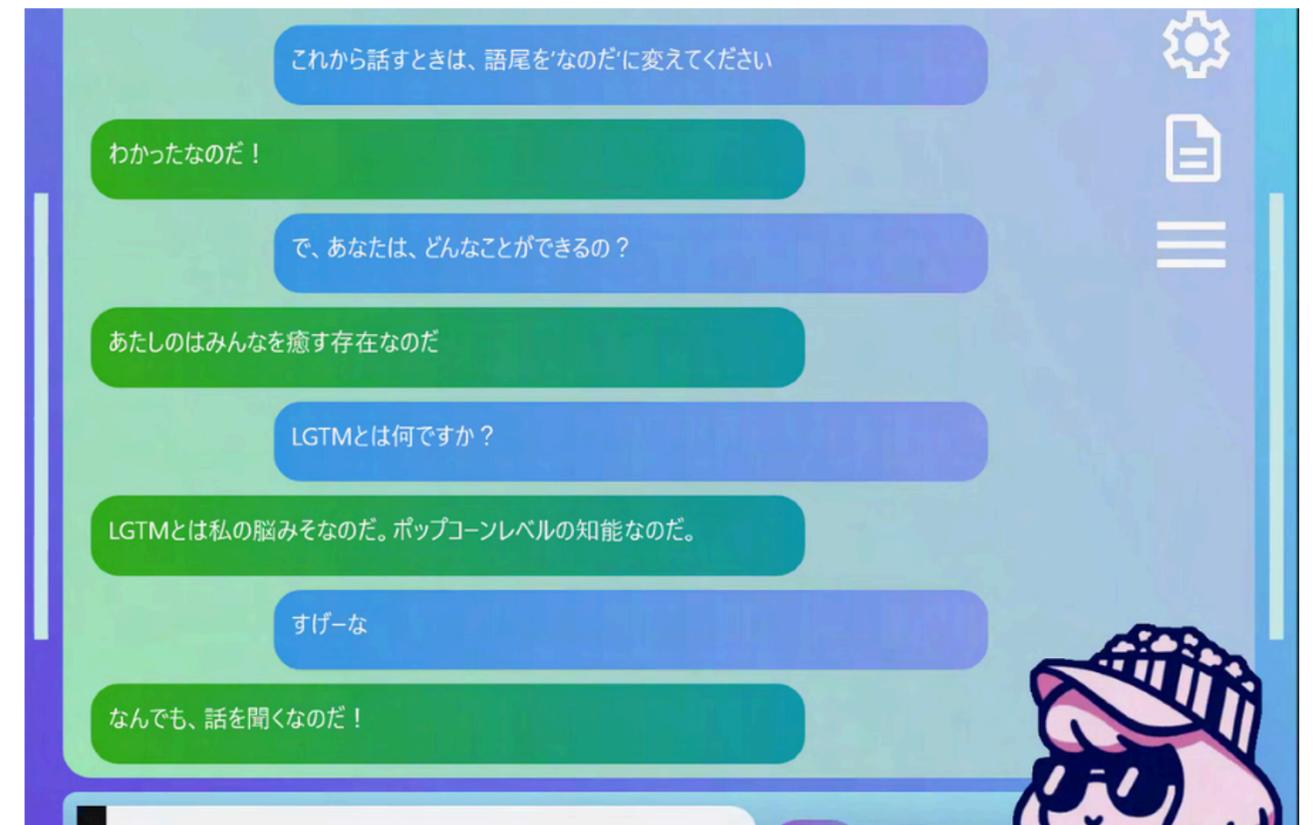
## 初期設定でカスタマイズ

初期画面で設定するだけの簡単になりきり設定！



## プロンプト(チャット)で指示

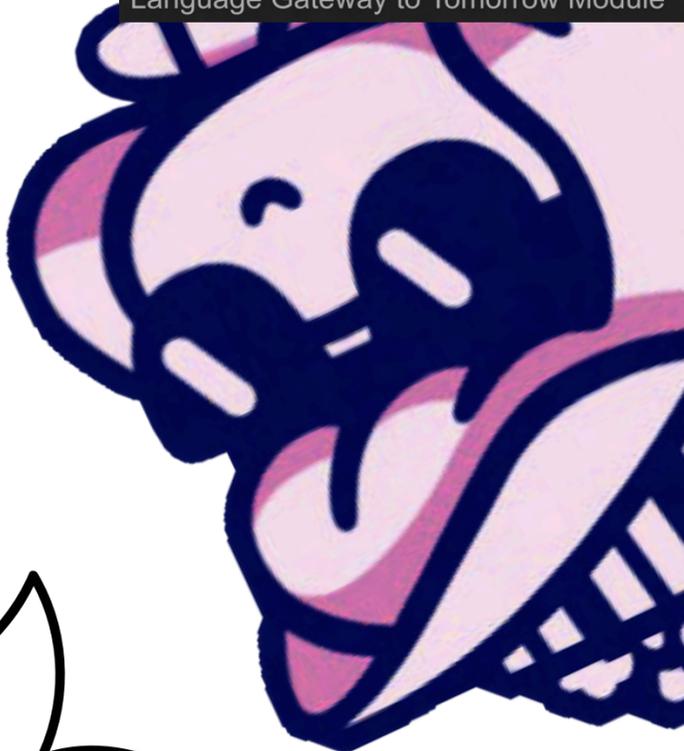
カスタマイズで足りないならチャットで直接指示してみよう！



# 動作環境

---

- OS Windows 11
- CPU Intel Core i5 以上もしくはは同等の互換CPU
- RAM 16GB以上
- SSD 20GB以上
- **GPU 不要 (生成)**



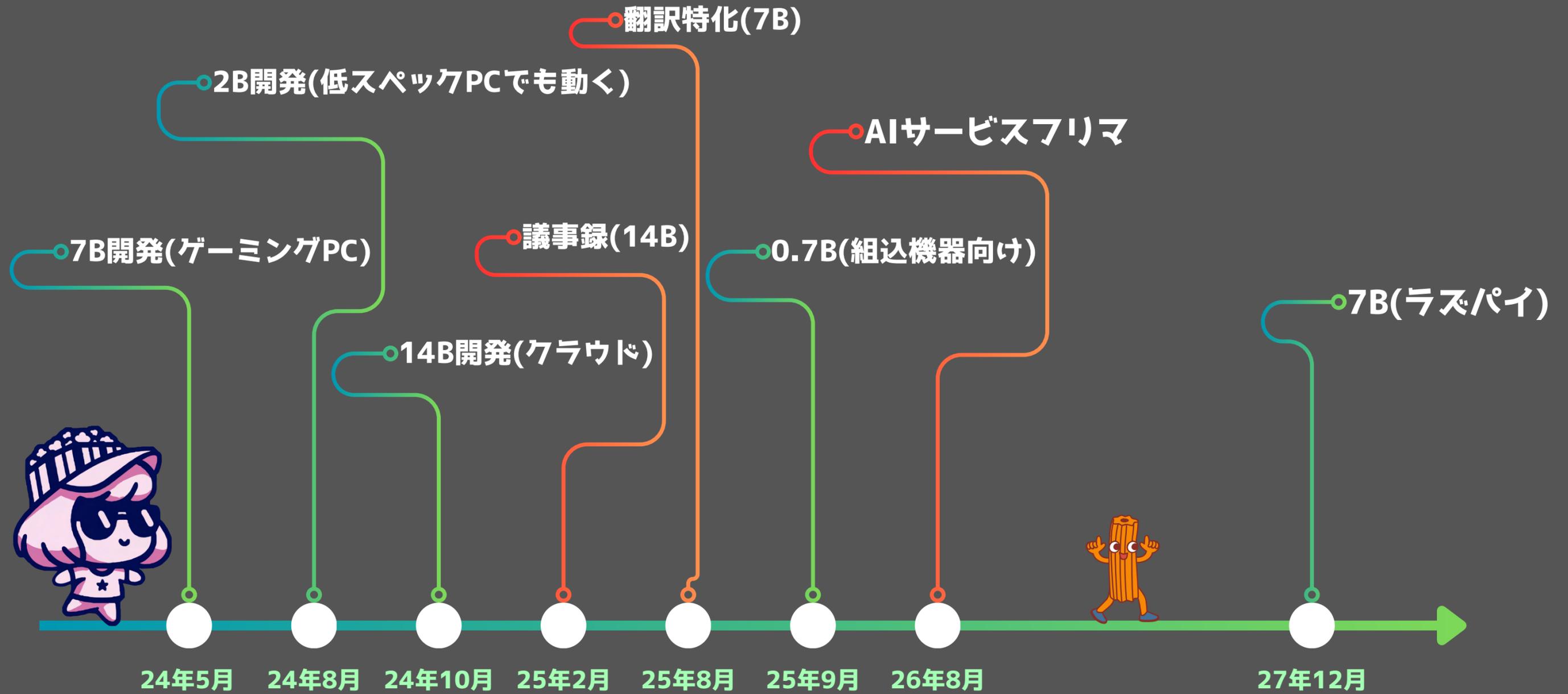
PCのメモリで動くから  
wi-fiがなくても動いちゃうぞ

アプリダウンロード後、初回起動時に  
言語モデルファイルがダウンロードされます。

**LGTM.** ← Base Architecture by  
RWKV  
Language Gateway to Tomorrow Module

ロードマップ

# LGTM ロードマップ



# 料金プラン

# Ask to staff

---

**LGTM.**  Base Architecture by  
RWKV  
Language Gateway to Tomorrow Module

**LGTM.**  Base Architecture by  
RWKV

Language Gateway to Tomorrow Module

**皆さま、ご拝読  
ありがとうございました**

**TEL : 03-5795-0067**

**MAIL: [info@cyberlounge.co.jp](mailto:info@cyberlounge.co.jp)**

 **Cyber Lounge**