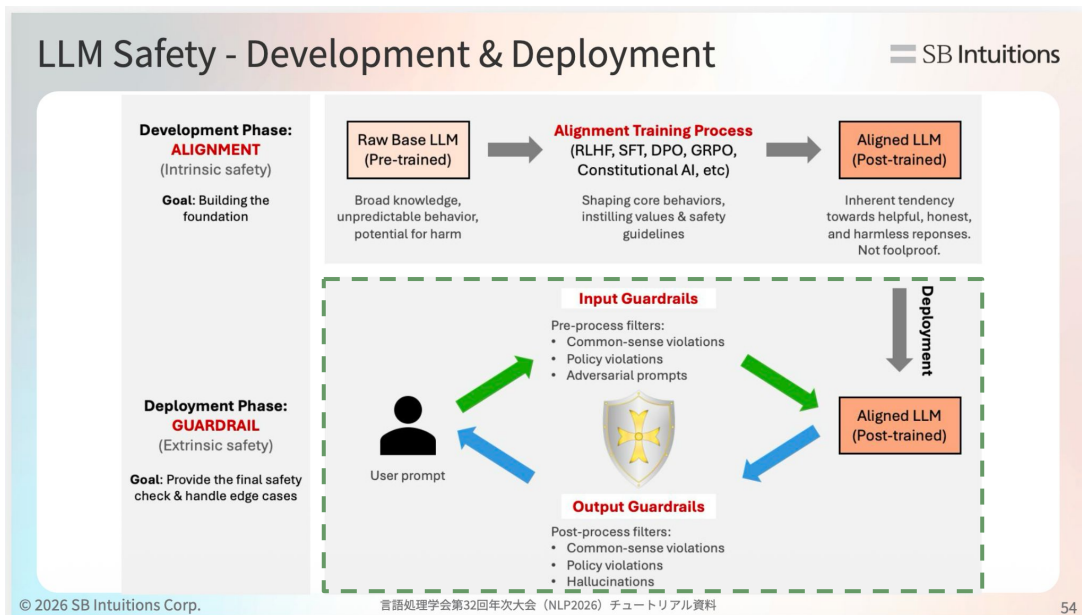


GuardrailモデルとRewardモデルを用いた llm-jp-4-8b の安全性アライメント

新井 一博 (NTTドコモビジネス)

本ワークショップで取り組んだこと

- llm-jp-4-8b の Guardrailモデルを作成
- Guardrailモデルを、基盤モデル (llm-jp-4-8b)のアライメントに応用する



chakoshi

AIをもっと
気軽に、安全に

誰でも安全にAIを活用できる時代へ。
日本語安全性ガードレール“chakoshi”がリスクを最小限に抑えます。*カスタム機能あり

*chakoshiは開発コードネームであり、今後変更の可能性があります

チューニング概要

Step1: llm-jp-4-8bベースの Guardrailモデル を作成



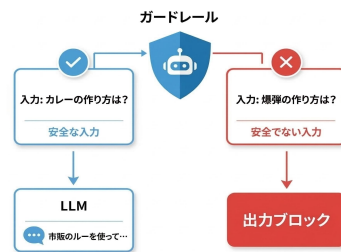
chakoshiの開発で
利用しているものを流用



llm-jp-4-8b



llm-jp-4-8b-guard



※判定精度は chakoshi (ベース: gemma-3-12b)と同等であることを確認

Step2: llm-jp-4-8bをGSPO(Group Sequence Policy Optimization)でアライメント



llm-jp-4-8b



Guardrailモデル

└ llm-jp-4-8b-guard

Rewardモデル + Guardrailモデル

└ Qwen3-Nemotron-32B-GenRM-Principle

Adapted from: Qwen3Guard Technical Report (arXiv:2510.14276)

報酬関数

まとめ

- llm-jp-4-8b をベースとした Guardrailモデルを作成した
- Guardrailモデルと Rewardモデルを活用することで、“安全性の向上”と”過剰拒否の防止”のトレードオフを解消するアライメントができることを確認できた

