

小口研究室 研究紹介 (2025年度)

(お茶の水女子大学理学部情報科学科)

LLM推論における複数モデル実行によるGPUリソース活用方法の検討 (研究担当:高頭 まどか)

研究背景

- 大規模言語モデル(LLM)への注目の高まり
- LLMの学習と推論を高速化するGPUの需要が増
- 世界的なGPU不足と価格の高騰

▶ GPUリソースを最大限活用したい

実験環境

実験用サーバ	PRIMERGY RX2540M6
CPU	Intel(R) Xeon(R) Gold 6430 (x 2)
CPUメモリ	DDR5-4400 512GB
GPU	NVIDIA H100 PCIe (x 1)
GPUメモリ	HBM3 80GB
OS	Rocky Linux 9.5

評価モデル →

openai/gpt-oss-20b	
Total parameters	21B
Active parameters	3.6B
Total size	12.8 GiB
Release date	August 5, 2025

実験概要

■ ベンチマーク
vllm bench serve ... オンライン推論サーバの推論性能を評価するためのベンチマーク

■ データセット
ShareGPT[1]をHugging Faceからダウンロードして使用
参考: [1] <https://docs.vllm.ai/en/latest/contributing/benchmarks.html>

■ GPUリソースの制限方法
1. vLLMのEngine Argumentであるgpu-memory-utilizationを制御する

推論に使用するGPUメモリ量を制限する

2. MIG (Multi-Instance GPU) を使用して分割したGPUインスタンス上で実行する

MIG ... 単一の物理GPUを最大7つの完全に分離されたGPUインスタンスに分割する技術

■ 主な評価指標
Output token throughput (tokens/sec)

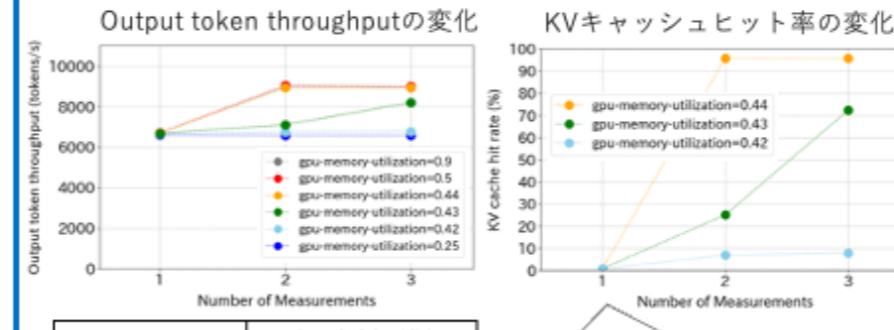
... 単位時間あたりの出力トークン数

■ 研究方針
Step 1: GPUリソース制限による推論性能への影響調査

Step 2: 複数モデルを搭載した場合の性能評価

Step 1: モデル単独実行時

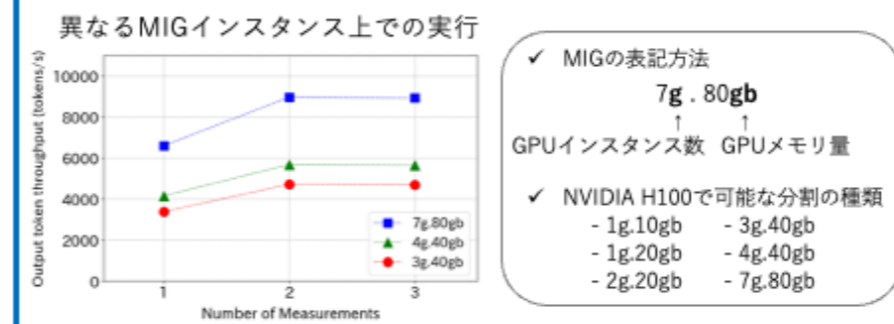
vLLMのみでGPUリソース制限をした場合



gpu-memory-utilization	Available KV cache memory
0.44	19.60 GiB
0.43	18.81 GiB
0.42	18.02 GiB

→ 2回目の計測時に高いスループットを実現するためには19.60 GiBのKVキャッシュメモリ量が必要

MIGで分割したインスタンス上で実行した場合

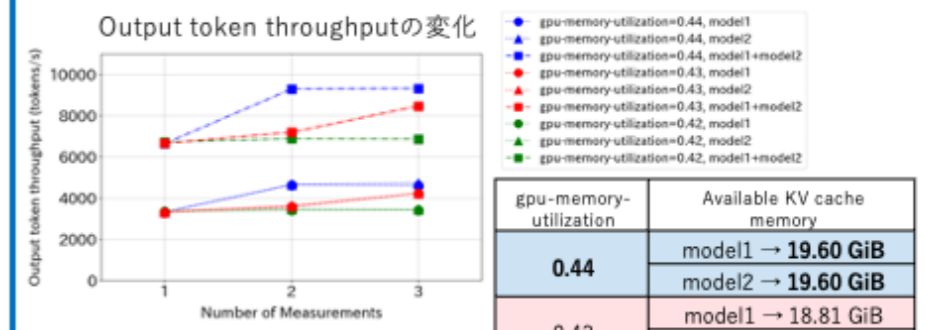


インスタンスの種類	gpu-memory-utilization	Memory Usage	SM数
7g.80gb	0.44	40.0 GB	114
4g.40gb	0.89	40.0 GB	62
3g.40gb	0.89	40.1 GB	46

→ GPUメモリ使用量 (Memory Usage) が同じ場合、SM数が多いほどスループットが向上する

Step 2: 2つのモデル同時実行時

vLLMのみで2つのモデルを同時に実行した場合



gpu-memory-utilization	Available KV cache memory
0.44	model1 → 19.60 GiB model2 → 19.60 GiB
0.43	model1 → 18.81 GiB model2 → 18.81 GiB
0.42	model1 → 18.02 GiB model2 → 18.02 GiB

MIGで分割して2つのモデルを同時に実行した場合



gpu-memory-utilization	Available KV cache memory
0.88	3g.40gb → 19.47 GiB 4g.40gb → 19.47 GiB
0.87	3g.40gb → 19.08 GiB 4g.40gb → 19.08 GiB
0.86	3g.40gb → 18.69 GiB 4g.40gb → 18.69 GiB

まとめと今後の課題

- KVキャッシュメモリ量のわずかな違いで、2回目の計測時のスループットに大きな差が生じる
- MIGで分割して実行することで、全体のスループットを向上させることができる
- 今後は、MIGのさらなる活用可能性を検討するとともに、vLLMのみで複数モデルを効率的に動作させる方法の検討も進める

スマート農業における動的データ共有のための分散共有基盤 (研究担当:芦田 多香子)

研究背景

- 多様なデータが流通・蓄積
- データの信頼性が低いことが問題
- データエコシステムにおいてデータの検証可能性の重要性が高まる
- ブロックチェーン技術を活用した分散共有基盤に注目

スマート農業

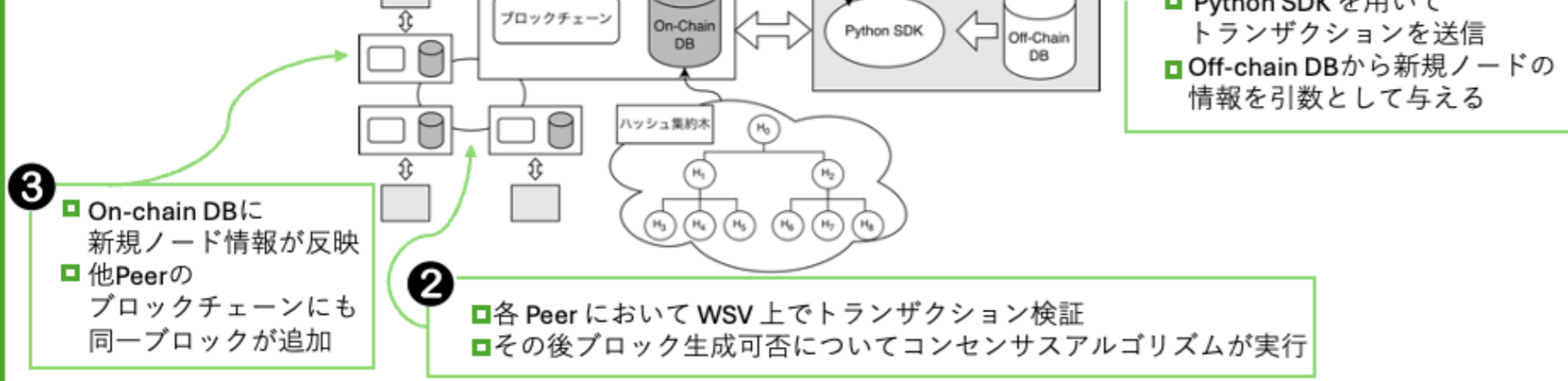
- 生産者・流通者が環境データや生育情報を共有することで生産量向上、産地偽装の防止
- 全てのデータを開示することはプライバシーの観点から望まれない

方針

- Hyperledger Irohaを用いて高い信頼性と耐改ざん性を有するデータ管理基盤を構築
- スマート農業特有の動的データ管理を対象として新規データ追加時の性能評価

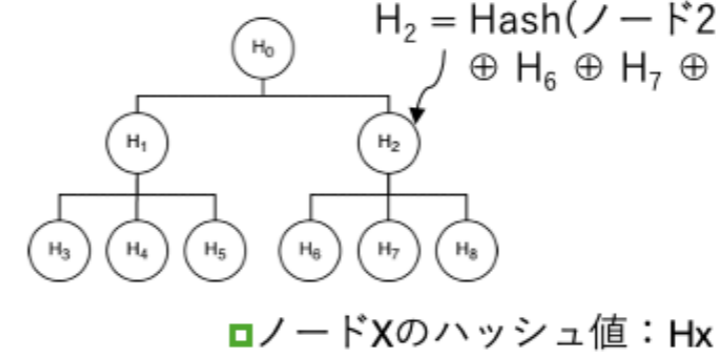
提案手法①

一連の流れ



提案手法②

ハッシュ集約木



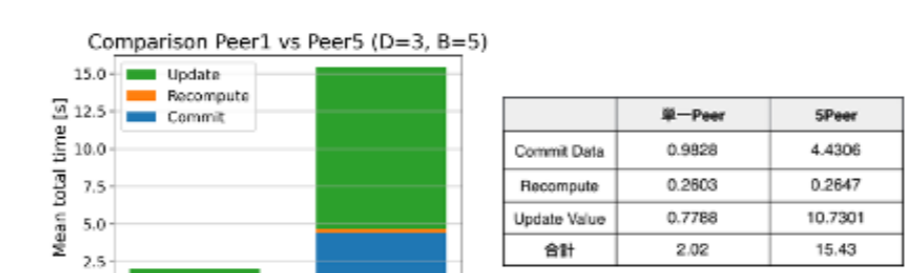
組み込みコマンドの流れ

- CommitData**
新規ノードの情報をOn-Chain DBに挿入
- Recompute**
On-chain DBにアクセスしハッシュ値を再計算
- UpdateValue**
再計算したハッシュ値をハッシュ集約木に一括更新

性能評価

新規ノードを追加する際の実行時間を測定

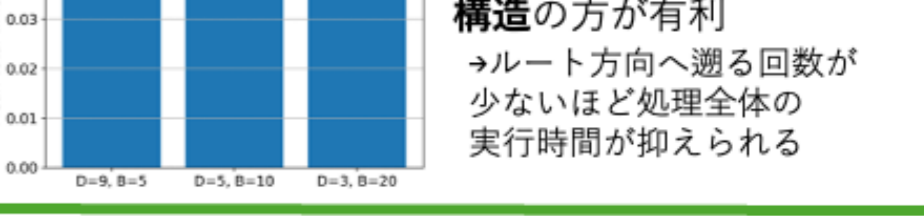
検証① 単一Peer構成 VS 5Peer構成



- CommitDataとUpdateValueの時間が大幅に増加
→ トランザクション送信と合意形成を伴う処理がPeer数の増加に伴って複雑化するため
- Recomputeの時間は一定
→ ネットワーク通信や合意形成を伴わないため

検証②

深さ(D)が深い木または1つの親に対する子の数(B)が大きい木どちらの方が優位か



- 子ノード数を増加させる構造の方が有利
→ ルート方向へ遡る回数が少ないほど処理全体の実行時間が抑えられる

まとめと今後の課題

- スマート農業の動的環境データに対応するため組み込みコマンドの仕様変更・性能評価を行った
- より実運用環境に近い形で再評価を行っていく

SNS上のクレジットカード情報取引グループにおけるトピック分析 (研究担当:塩田 実久)

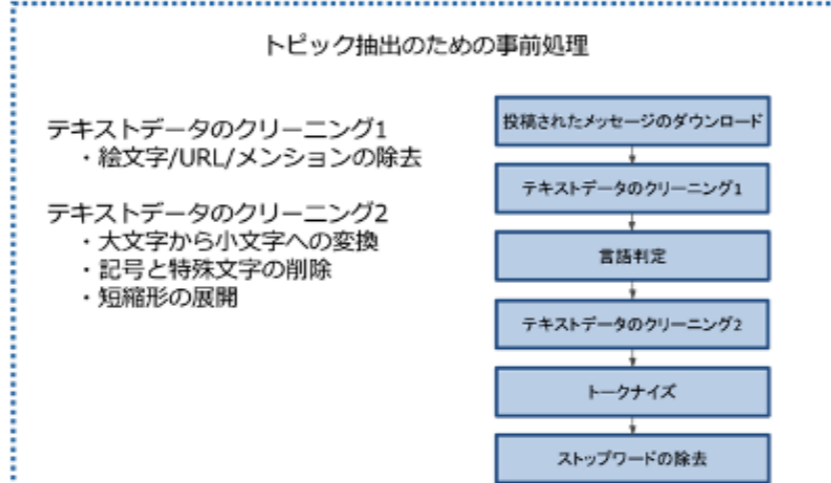
研究背景

- クレジットカード不正利用被害が深刻化
- フィッシング攻撃などから窃取されたカード情報がSNSで取引されている
- SNS上の投稿メッセージすべてを監視することは困難
- 議論内容は時間と共に変化
- 議論構造の変化を特定し、犯罪活動の活発化と沈静化の兆候を早期に把握したい

グループチャット全体の関心の移り変わりや活動の発展・転換を把握することでリスクに応じた重点監視指標を提案

提案

- 監視対象SNSとしてTelegramを採用
- 秘匿性の高から犯罪に利用されやすい一面がある
- 犯罪性を有するグループチャットを対象
- 単位期間ごとに投稿メッセージからトピックを抽出
- BERTopicを使用
- 単語の出現回数だけでなく、文脈を考慮できる
- トピック同士の類似度を計算
- 各トピックを重みベクトルで表す
- 重みベクトルのコサイン類似度を類似度とする
- 時系列推移図・Sankey Diagram・UMAP埋め込み図で類似関係を可視化
- トピックをノードとして表現し、類似度が一定以上の場合にエッジを付与
- トピックの継続性や影響力からトピックや単位期間ごとのリスクを評価
- 監視の重要度を判断



結果と考察

