# Private Cloud Computing

Jiale Guan

2024-08-12

# Outline

1. Privacy Leakage Surface
   - Architecture

2. Private Cloud Compute
   - Taxonomy
   - Requirements

3. LLM Serving Systems
   - Optimization Goals
   - Violations
   - Optimizations

# Outline

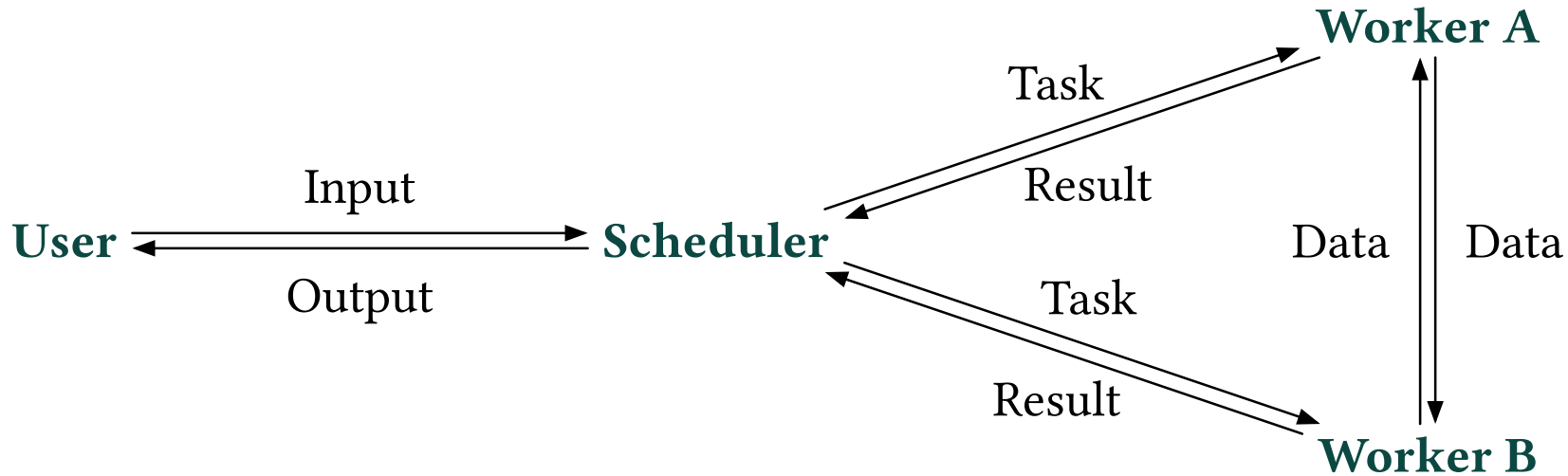1. Privacy Leakage Surface
   - Architecture

2. Private Cloud Compute
   - Taxonomy
   - Requirements

3. LLM Serving Systems
   - Optimization Goals
   - Violations
   - Optimizations

## 1.1 Architecture

# Outline

# 2.1 Taxonomy

| Requirements | Threat | Guarantees |
|---|---|---|
| Stateless computation | Trace of data after processing<br>e.g. Logging, debugging | (Purpose) Only use user data to perform requested operations<br>(Transient) Delete the data after fulfilling the request<br>(Scope) Not available to even Apple staff |
| Enforceable guarantees | Technical enforceability<br>e.g. External TLS-terminating load balancer | (Hardware) Secure Enclave, Secure Boot<br>(System) Signed System Volume, Swift on Server<br>(Software) Code Signing, Sandboxing |
| No privileged runtime access | Privileged interfaces<br>e.g. Shell access by SREs | No remote shell. Only pre-specified, structured, and audited<br>logs/metrics can leave the node<br>User data is reviewed by multiple indepedent layers |
| Non-targetability | Targeted attack<br>e.g. Steer request to compromised nodes | (Hardware) Hardened supply chain<br>(Scheduler) Requests cannot be user/content-specific routed<br>(Anonymity) OHTTP Relay, RSA Blind Signature<br>(Scope) No system-wide encryption |
| Verifiable transparency | Uninspected code | Every production build of PCC publicly available |

# 2.2 Requirements

## Stateless computation

Private Cloud Compute must use the personal user data that it receives exclusively for the purpose of fulfilling the user's request. This data must never be available to anyone other than the user, not even to Apple staff, not even during active processing. And **this data must not be retained**, including via logging or for debugging, after the response is returned to the user. In other words, we want a strong form of stateless data processing where **personal data leaves no trace** in the PCC system.

# Enforceable guarantees

Security and privacy guarantees are strongest when they are entirely technically enforceable, which means it must be possible to **constrain and analyze all the components** that critically contribute to the guarantees of the overall Private Cloud Compute system. To use our example from earlier, it's very difficult to reason about what a TLS-terminating load balancer may do with user data during a debugging session. Therefore, PCC must not depend on such external components for its core security and privacy guarantees. Similarly, operational requirements such as collecting server metrics and error logs must be supported with mechanisms that do not undermine privacy protections.

# No privileged runtime access

Private Cloud Compute must not contain privileged interfaces that would enable Apple's site reliability staff to bypass PCC privacy guarantees, even when working to resolve an outage or other severe incident. This also means that PCC must not support a mechanism by which the privileged access envelope could be enlarged at runtime, such as by loading additional software.

# Non-targetability

An attacker should not be able to attempt to compromise personal data that belongs to specific, targeted Private Cloud Compute users without attempting a broad compromise of the entire PCC system. This must hold true even for exceptionally sophisticated attackers who can attempt physical attacks on PCC nodes in the supply chain or attempt to obtain malicious access to PCC data centers. In other words, a limited PCC compromise must not allow the attacker to **steer requests from specific users to compromised nodes**; targeting users should require a wide attack that's likely to be detected. To understand this more intuitively, contrast it with a traditional cloud service design where every application server is provisioned with database credentials for the entire application database, so a compromise of a single application server is sufficient to access any user's data, even if that user doesn't have any active sessions with the compromised application server.

# Verifiable transparency

Security researchers need to be able to verify, with a high degree of confidence, that our privacy and security guarantees for Private Cloud Compute match our public promises. We already have an earlier requirement for our guarantees to be enforceable. Hypothetically, then, if security researchers had sufficient access to the system, they would be able to verify the guarantes. But this last requirement, verifiable transparency, goes one step further and does away with the hypothetical: security researchers must be able to verify the security and privacy guarantees of Private Cloud Compute, and they must be able to verify that the software that's running in the PCC production environment is the same as the software they inspected when verifying the guarantees.

# Outline

# 3.1 Optimization Goals

- Lower resource requirements
  - ▸ Model compression
  - ▸ Collaborative inference
  - ▸ Offloading to CPU and disk
- Higher throughput
  - ▸ Batch
  - ▸ Tensor parallelism
  - ▸ Sequence parallelism (Input sequence)
  - ▸ Pipeline parallelism (Stage placement)
- Lower latency
  - ▸

# 3.2 Violations

| Requirement | Violations |
|---|---|
| Stateless computation | Logging, prioritization, history metadata |
| Enforceable guarantees | Data transfer/duplication, data offloading, access control |
| No privileged runtime access | Monitoring, debugging, profiling |
| Non-targetability | Biased scheduler, input/output leakage |
| Verifiable transparency | Uninspected code |

Universal problems: Access control between worker nodes.

| Requirement | Violations | FT 22 | FlexGen 23 | Sarathi 23 | Mooncake 24 |
|---|---|---|---|---|---|
| Stateless computation | Logging | | | | |
| | History metadata | | | | |
| Enforceable guarantees | Data transfer/duplication | | | | ✓ |
| | Data offloading | | ✓ | | ✓ |
| No privileged runtime access | Monitoring | | | | |
| | Debugging | | | | |
| | Offline Profiling | | | | ✓ |
| Non-targetability | Priority-based scheduler | | | | ✓ length |
| | Input/output leakage | | | | |
| Verifiable transparency | Uninspected code | | | | ✓ |

# 3.3 Optimizations

**S**: Stateless computation **E**: Enforceable guarantees **P**: No privileged runtime access **T**: Non-targetability **V**: Verifiable transparency

| Type | Optimization | Threat | FT 22 | Orca 22 | vLLM 23 | FlexGen 23 | FastServe 23 | Sarathi 23 | DistServe 24 | TetriInfer 24 | Mooncake Moonshot | DeepSpeed Microsoft | TensorRT-LLM NVIDIA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Storage** (KVCache) | Paging | | | | Initial | | | ✓ | | | ✓ | ✓ | |
| | Offloading | SE | | | | ✓ | ✓ | | | | ✓ | | ✓ |
| | Duplication | T | | | | | | | | | ✓ | | |
| | Pulling | SET | | | | | | | ✓ | | | | |
| **Scheduler** | Priority-Based | T | | | | | ✓ | ✓ | | ✓ | ✓ | ✓ | |
| | Local Scheduler | ET | | | | | ✓ | | | ✓ | ✓ | | |
| | Instance Flip | P | | | | | | | | ✓ | | | |
| | Disaggregated Arch | | | | | | | | ✓ | ✓ | ✓ | | |
| | Online/Offline Profiling | P | | | | | | | ✓ | | ✓ | | |
| **PP** | Iteration-Level Batch | | | Initial | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| | Chunked Prefill | | | | | | | Initial | | ✓ | ✓ | ✓ | |
| | Prepack Prefill | | | | | | | | ✓ | ✓ | | | |
| **Verification** | Uninspected Code | V | | | | | | | | | ✓ | | |

# Paging & Offloading

Definition:
- Paging: 使用类似于虚拟内存的机制，将模型参数分页存储在磁盘上，根据需要加载到内存中。
- Offloading: 将模型参数从 GPU 内存中移动到 CPU 内存或磁盘中，以释放 GPU 内存供其他模型使用。

Threats:
- 分页处理过程中可能会产生包含敏感信息的日志，这些日志如果没有妥善管理，可能会泄露隐私数据。
- 分页数据可能会被意外持久化到不安全的存储介质中，从而暴露隐私数据。

# Duplication & Pulling

Definition:
- Duplication: 在不同的节点之间复制模型参数，以便在多个节点上并行执行推理任务。
- Pulling: 从远程节点拉取模型参数，以便在本地节点上执行推理任务。

Threats:
- 模型参数的复制和拉取过程中可能会泄露隐私数据。
- 模型参数的复制和拉取过程中可能会定向到恶意节点，从而导致隐私数据泄露。如果其中任何一个节点被攻破，攻击者可能获得整个模型的敏感信息。
- 拉取模型参数可能导致数据不同步，尤其在多次拉取操作之间，可能出现数据不一致的情况，影响模型的准确性和隐私保护。

# Priority-based Scheduler & Local Scheduler & Instance Flip

Definition:
- Priority-based Scheduler: 根据任务的优先级调度任务，以确保高优先级任务能够及时完成。
- Local Scheduler: 在本地节点上调度任务，以减少任务调度的延迟。

Threats:（优先级调度）
- 可能通过观察任务的优先级来推断任务的重要性和敏感性，从而有针对性地进行攻击。
- 在任务调度过程中，任务的调度信息（如任务类型、数据类型等）可能被泄露，导致隐私数据暴露。'

（本地调度）
- 在本地节点上调度任务时，所有任务和数据都集中在本地节点，如果本地节点被攻破，所有数据和任务信息都可能被泄露。
- 本地节点可能会缓存大量的任务数据，如果这些缓存数据未妥善处理，可能会导致隐私泄露。
- 为了减少调度延迟，可能会牺牲一些数据同步和一致性机制，导致数据不一致。

（节点翻转）
- 攻击者可能修改恶意节点的数据，来让恶意节点被选中执行任务，从而获取敏感信息。
- 攻击者可能通过控制节点翻转的时机，来获取敏感信息。

# Disaggregated Architecture & Online/Offline Profiling

Definition:

- Disaggregated Architecture: 将 Prefill 和 Decode 的过程通过实例（instance）分离，以提高资源利用率和灵活性。
- Online/Offline Profiling: 在线/离线性能分析，以优化模型推理性能。

Threats:

- 在进行用户画像时，会收集和存储大量的用户数据，包括在线行为数据和离线数据，这些数据一旦被泄露，可能对用户隐私造成严重威胁。

# Iteration-Level Batch & Chunked Prefill & Prepack Prefill

Definition:

- Iteration-Level Batch: 在迭代级别上进行批处理，以提高模型推理性能。
- Chunked Prefill: 将 Prefill 过程分块，以减少 Prefill 的延迟。
- Prepack Prefill: 预先打包 Prefill 数据，以减少 Prefill 的延迟。

Threats:

- N/A.

# Production LLMs

TorchServe - PyTorch

Triton HuggingFace TGI FastServe

Prompt Cache SGLang AttentionStore Preble

Pollux

**Refer to DistServe Related Work**

Thanks