

Morsels: Explicit Virtual Memory Objects

DIMES '23

Alexander Halbuer

Christian Dietrich

Florian Rommel

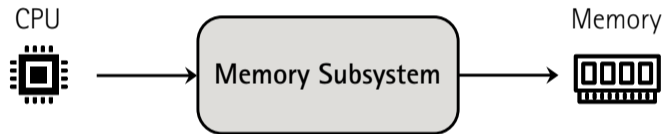
Daniel Lohmann

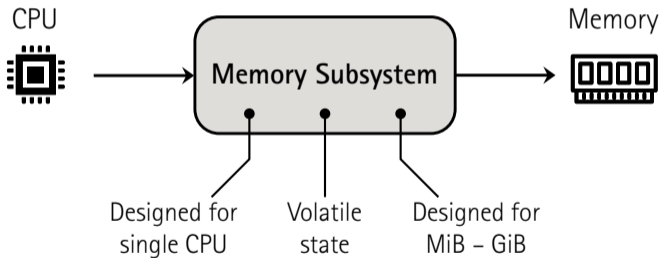


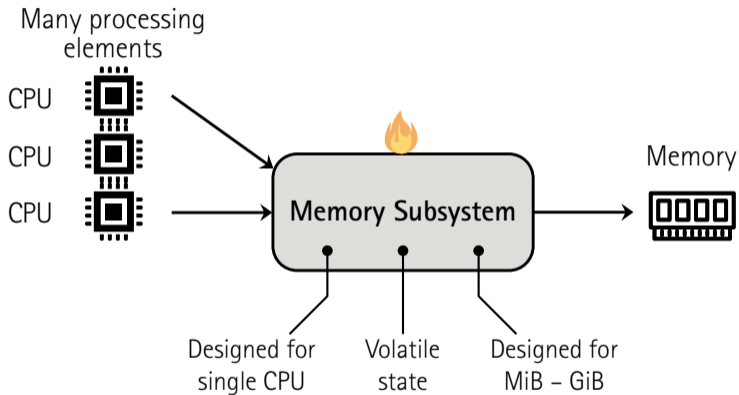
Leibniz Universität Hannover
halbuer@sra.uni-hannover.de

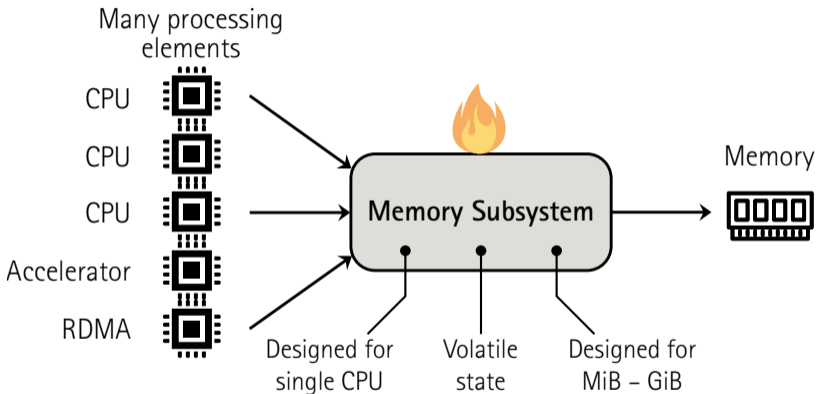
October 23, 2023

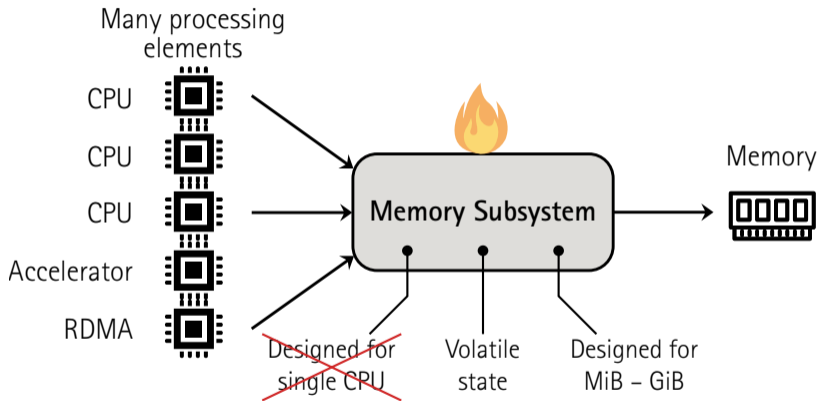


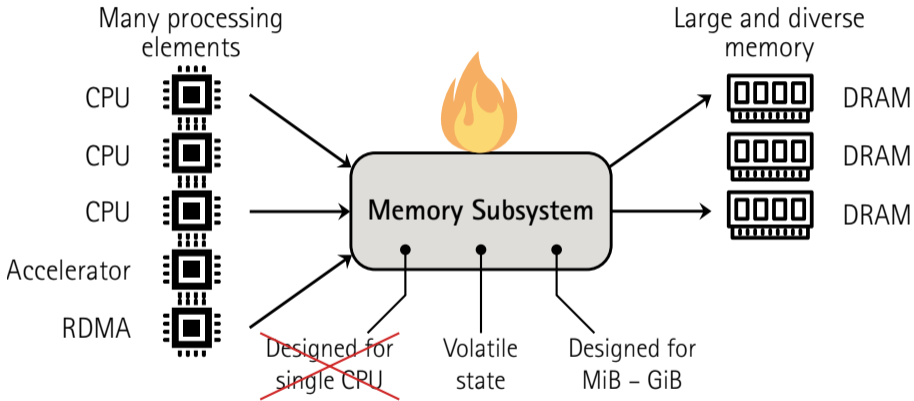


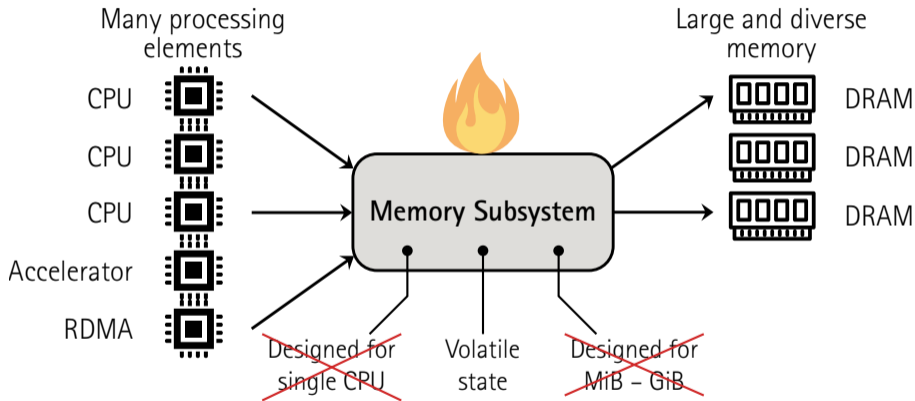


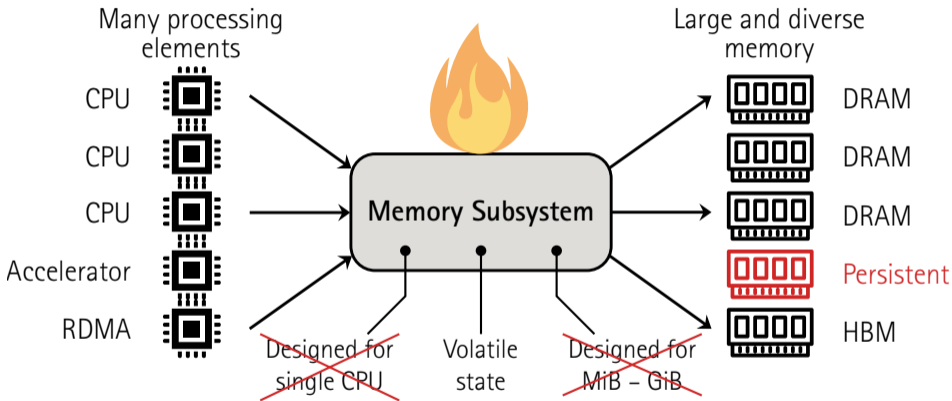


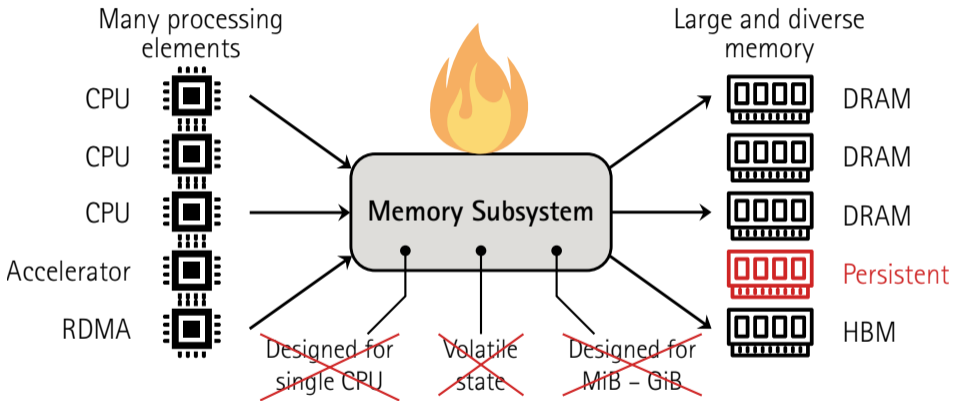


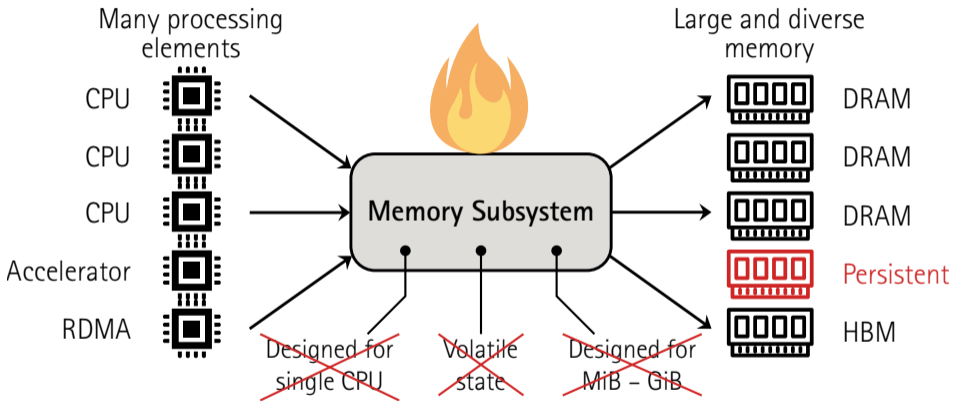




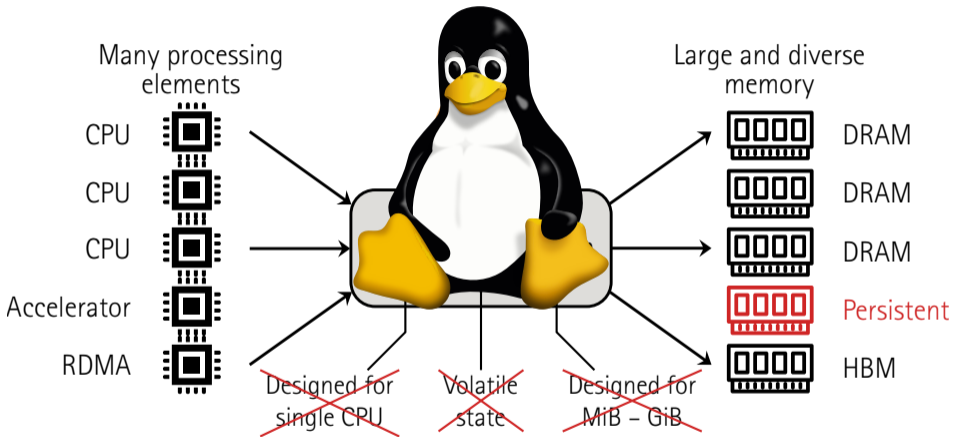








What are efficient OS abstractions for large memory objects?

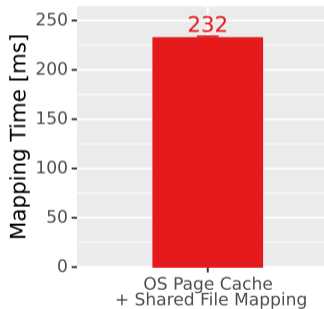


What are efficient OS abstractions for large memory objects?

Goal: Efficient handling of virtual memory

Goal: Efficient handling of virtual memory

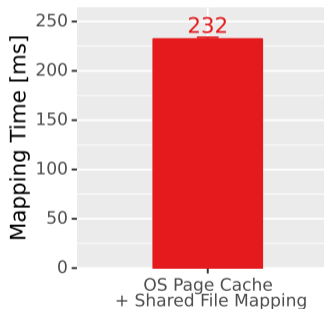
Mapping ~ 7 GiB of data into
an application's address space



Goal: Efficient handling of virtual memory

- Conventional file mapping (pre-faulted)
 - High setup costs
 - Create new page tables
 - Populate with pages from page cache

Mapping ~ 7 GiB of data into an application's address space

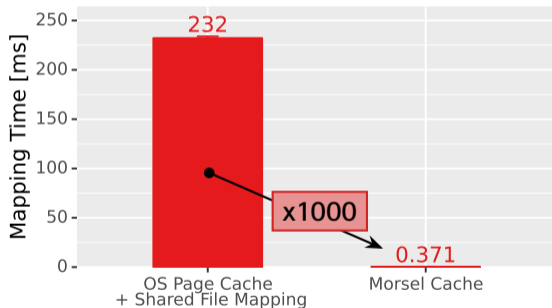


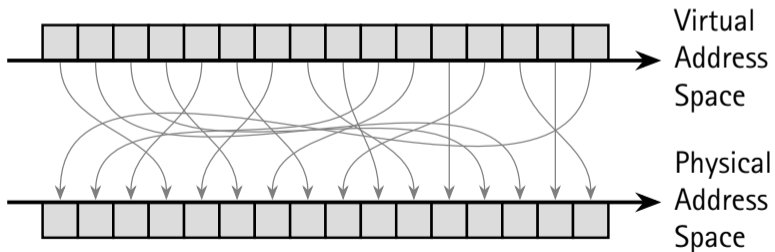
Goal: Efficient handling of virtual memory

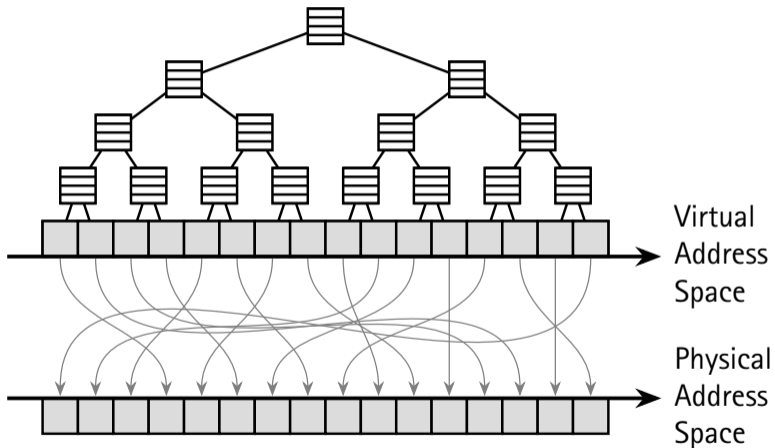
- Conventional file mapping (pre-faulted)
 - High setup costs
 - Create new page tables
 - Populate with pages from page cache

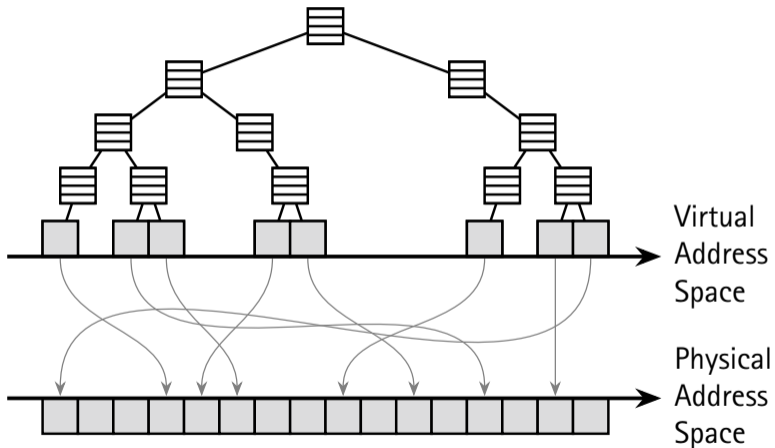
Our solution: Morsels

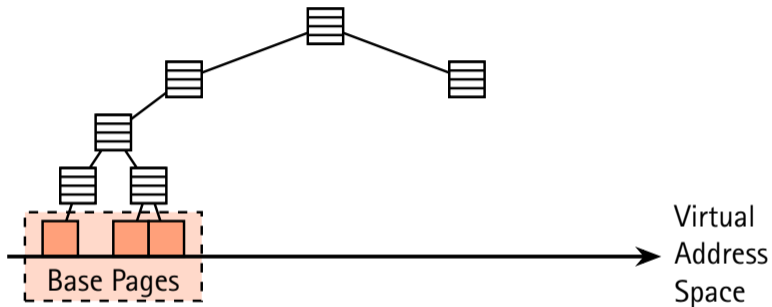
Mapping ~ 7 GiB of data into an application's address space

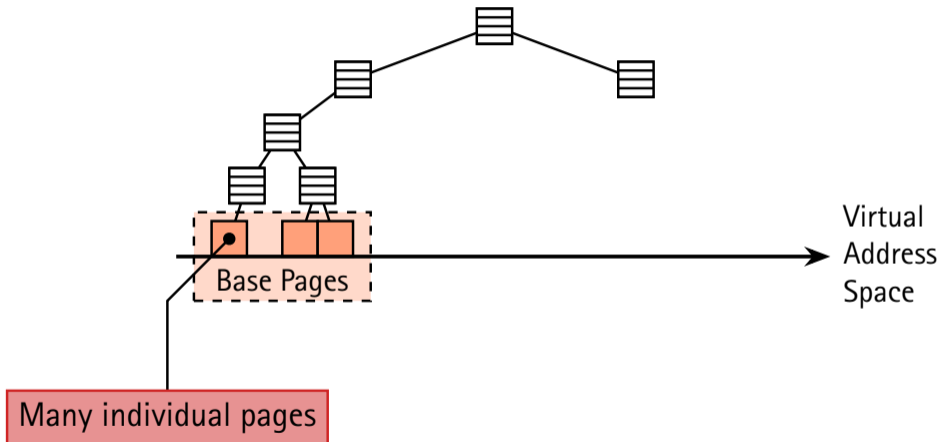


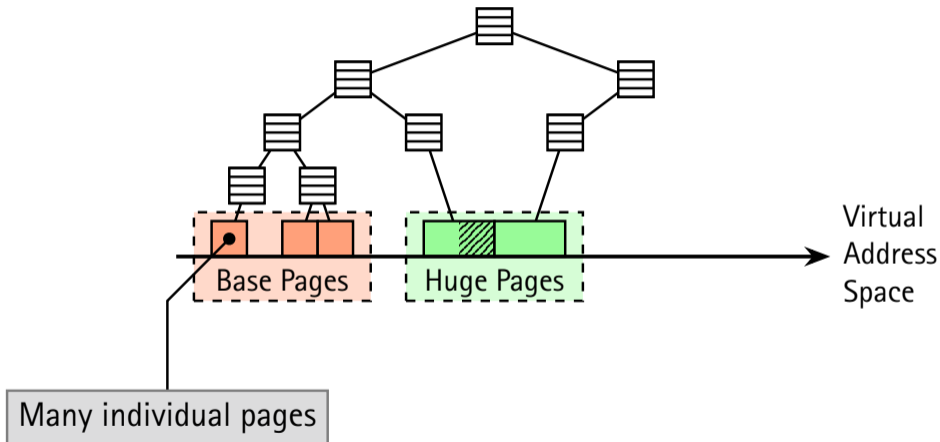


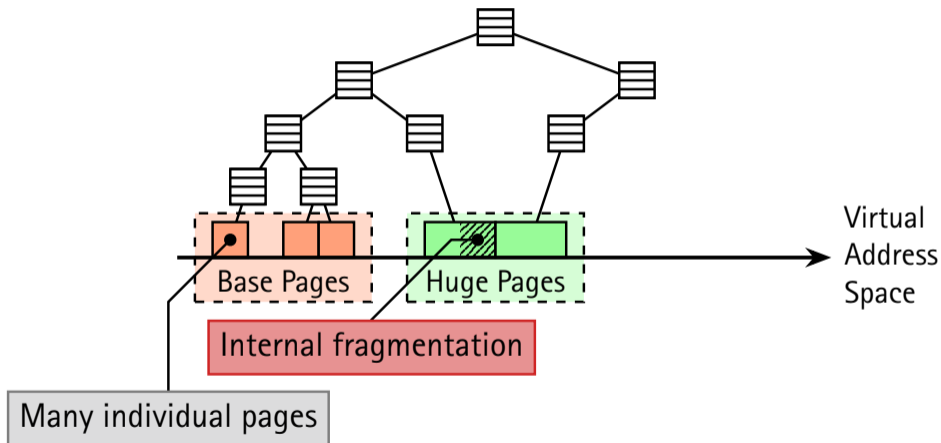


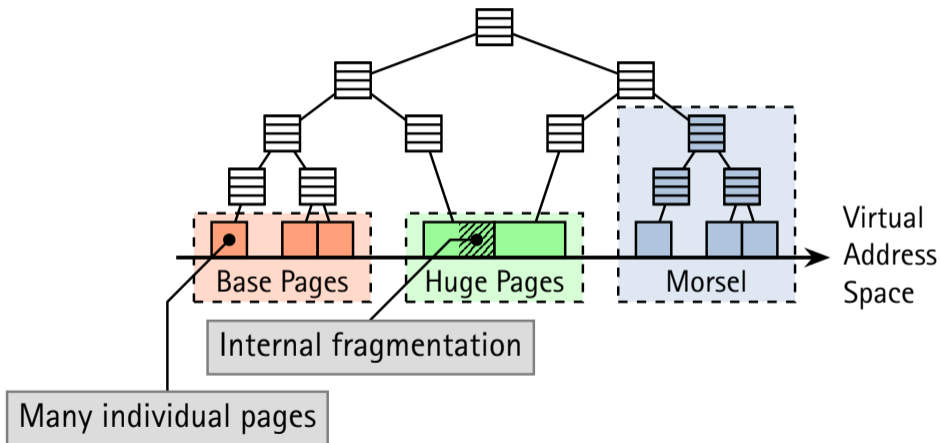


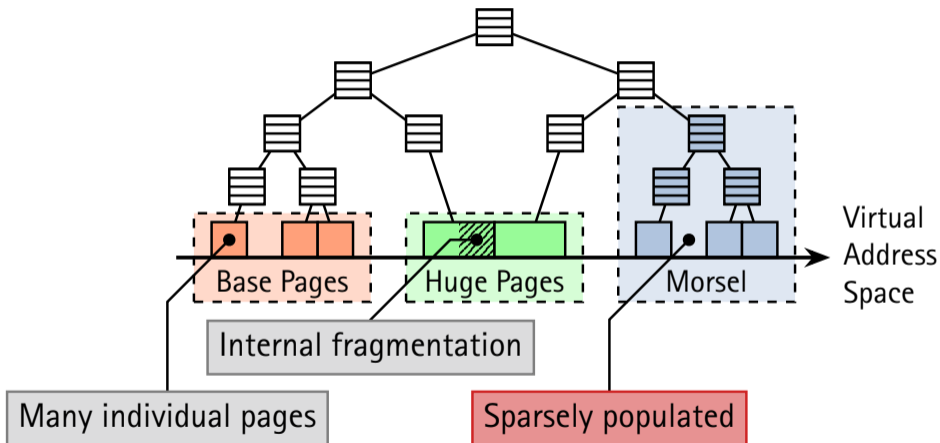


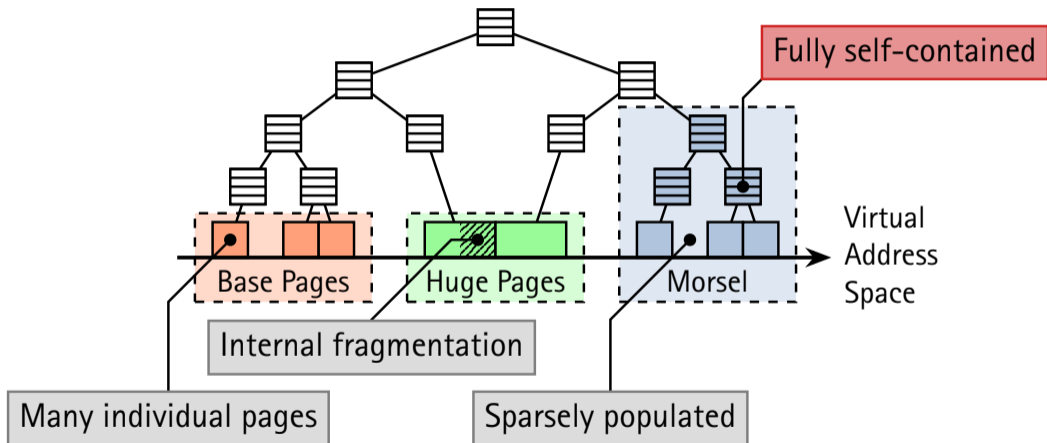


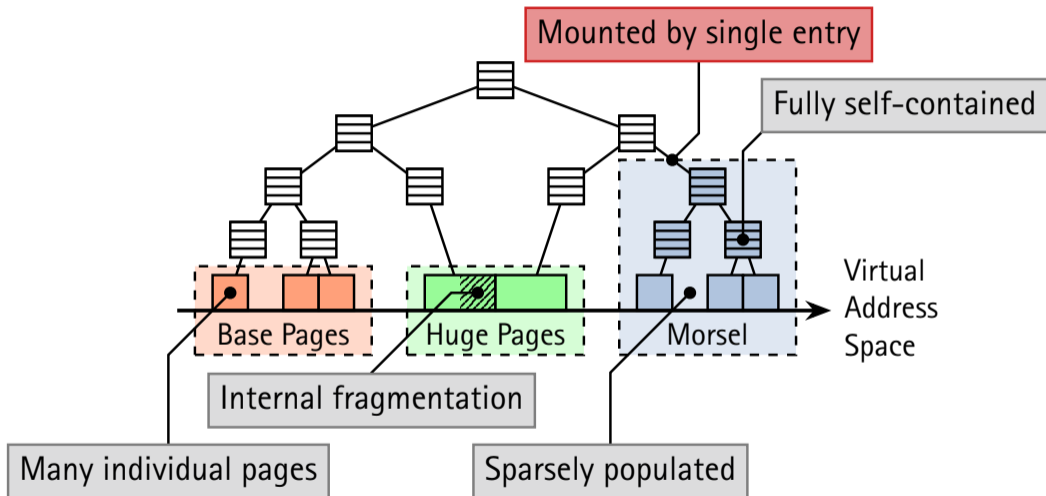


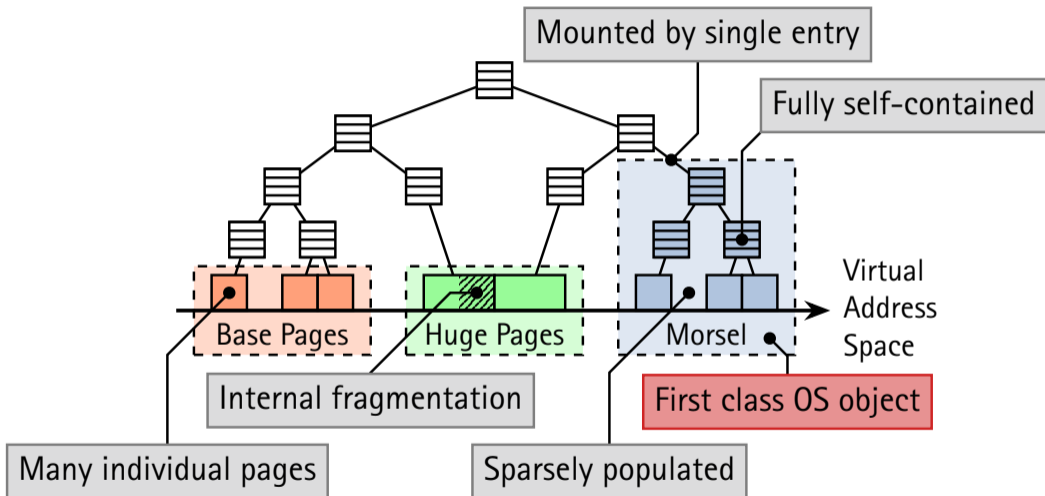


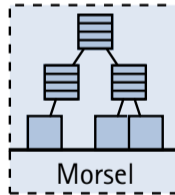




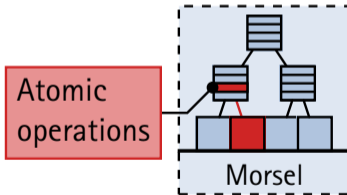




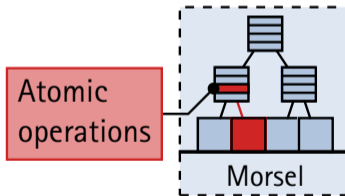




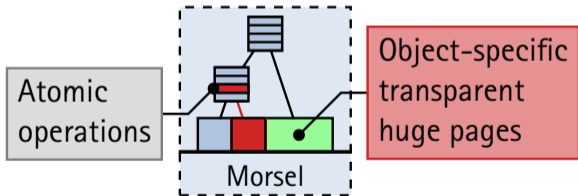
- Good multi-core scalability
 - No locking required



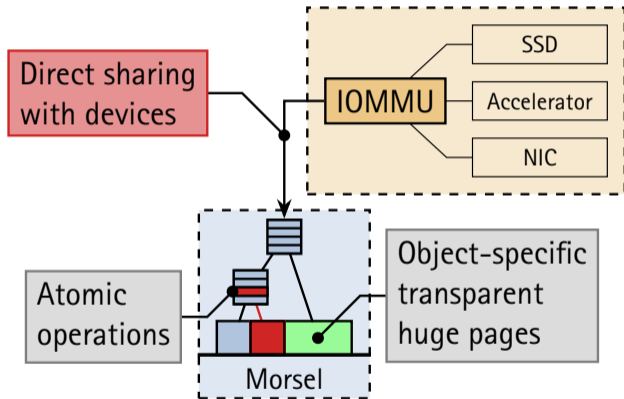
- Good multi-core scalability
 - No locking required
- Persistent objects on NVM
 - Fully self-contained design
 - Crash consistent implementation



- Good multi-core scalability
 - No locking required
- Persistent objects on NVM
 - Fully self-contained design
 - Crash consistent implementation
- Huge-page support
 - Increase TLB coverage
 - Transparency on object granularity



- Good multi-core scalability
 - No locking required
- Persistent objects on NVM
 - Fully self-contained design
 - Crash consistent implementation
- Huge-page support
 - Increase TLB coverage
 - Transparency on object granularity
- Sharing with devices
 - Reduce management overhead for DMA



Evaluation: Prototypic implementation in Linux

Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

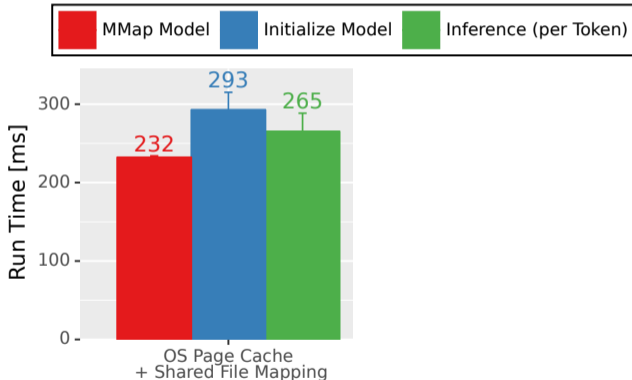
- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size

Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size

llama.cpp with SelfEe model: 6.82 GiB

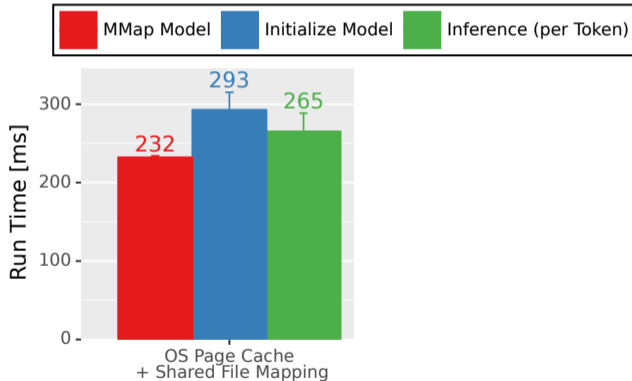


Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size
- Morsel mapping
 - Cache provides model data as Morsel
 - Single indivisible unit
 - Mapped in constant time

llama.cpp with SelfEe model: 6.82 GiB

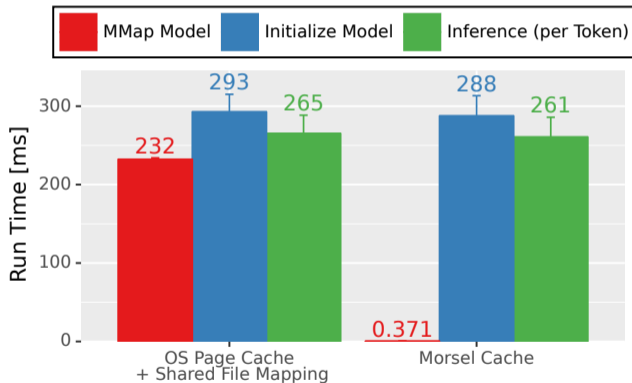


Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size
- Morsel mapping
 - Cache provides model data as Morsel
 - Single indivisible unit
 - Mapped in constant time

llama.cpp with SelfEE model: 6.82 GiB

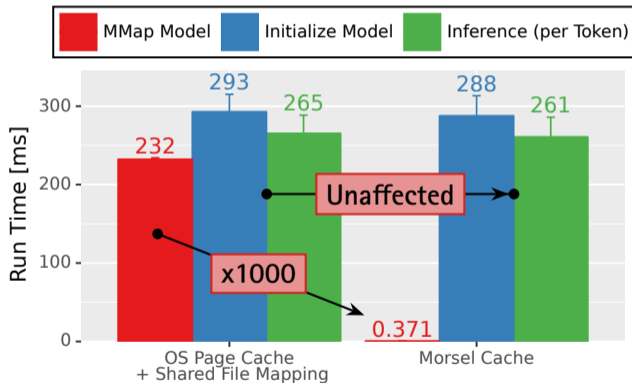


Evaluation: Prototypic implementation in Linux

Task: Inference with a large language model

- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size
- Morsel mapping
 - Cache provides model data as Morsel
 - Single indivisible unit
 - Mapped in constant time

llama.cpp with SelfEe model: 6.82 GiB



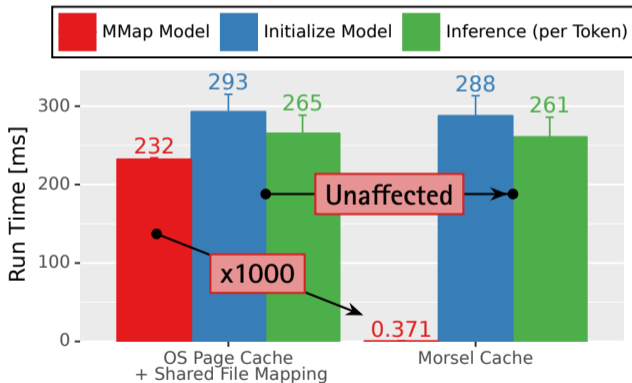
Evaluation: Prototypic implementation in Linux

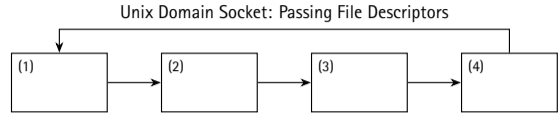
Task: Inference with a large language model

- Conventional file mapping (pre-faulted)
 - Uses individual 4-KiB pages
 - Populated from the page cache
 - Time increases with mapping size
- Morsel mapping
 - Cache provides model data as Morsel
 - Single indivisible unit
 - Mapped in constant time

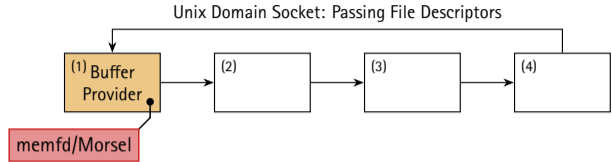
Result: 45% reduced startup time

llama.cpp with SelfEe model: 6.82 GiB

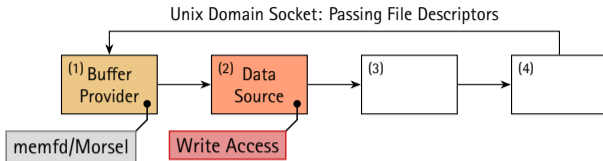




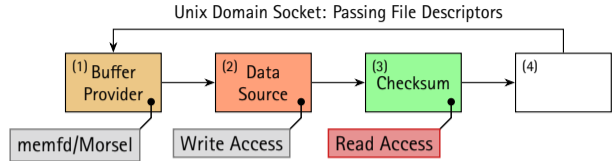
- Comparing Morsels with memfd



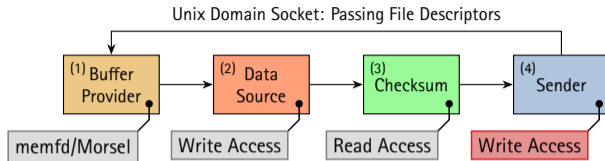
- Comparing Morsels with memfd



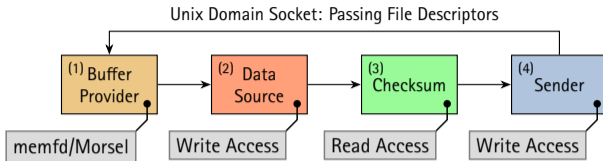
- Comparing Morsels with memfd



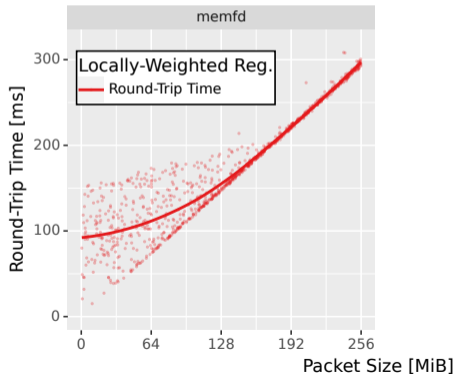
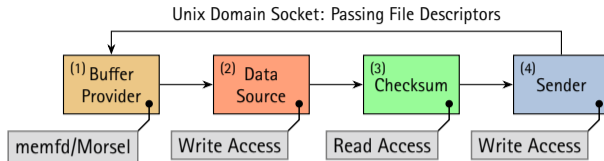
- Comparing Morsels with memfd



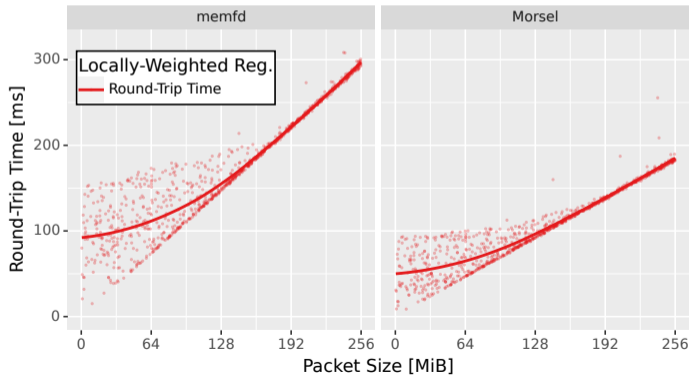
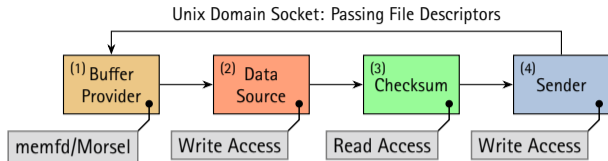
- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size



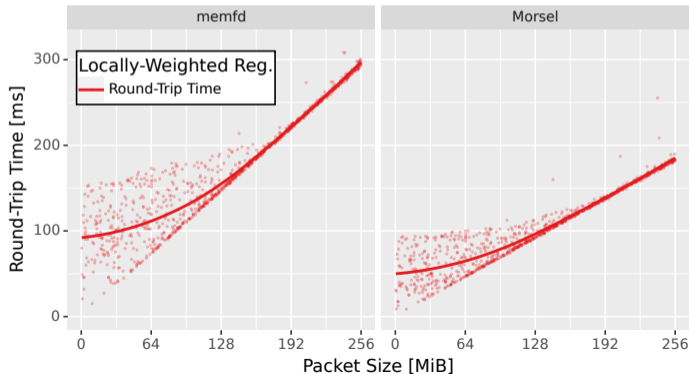
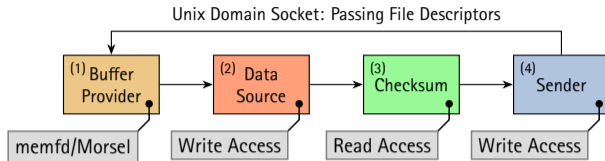
- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size



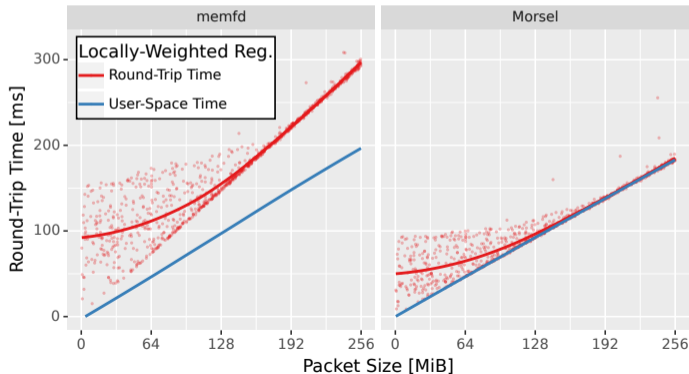
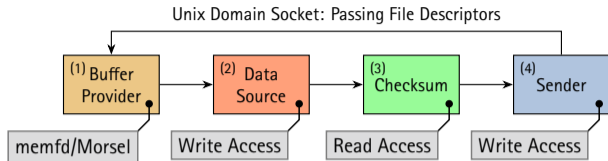
- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size



- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size
- Morsels: 39% reduced round-trip time on average



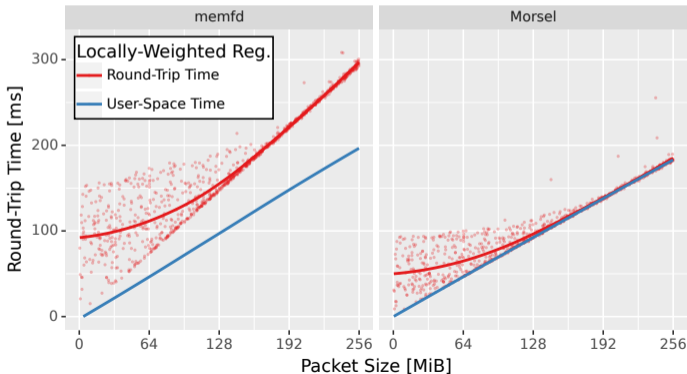
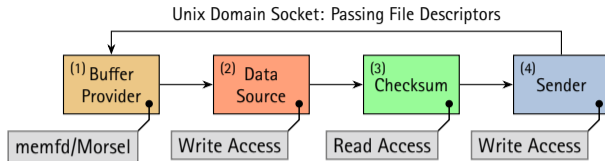
- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size
- Morsels: 39% reduced round-trip time on average



- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size

- Morsels: 39% reduced round-trip time on average

- Use-space time draws lower bound
 - + Waiting time
 - + Kernel processing
 - = Round-trip time

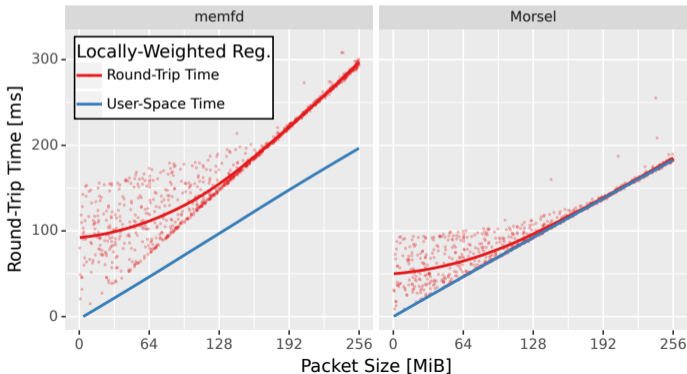
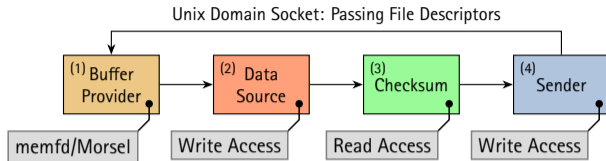


- Comparing Morsels with memfd
 - 1 000 packets per type
 - Up to 256 MiB packet size

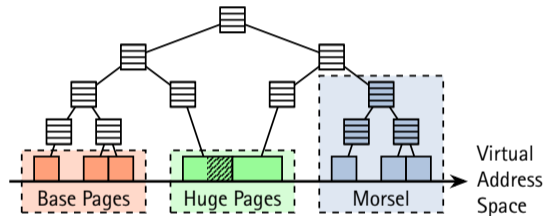
- Morsels: 39% reduced round-trip time on average

- Use-space time draws lower bound
 - + Waiting time
 - + Kernel processing
 - = Round-trip time

Conclusion: Morsels enable efficient handling of virtual memory!



- **Problem:** Conventional memory management
 - Many individual entities (pages)
 - Huge pages induce fragmentation
 - No persistent state

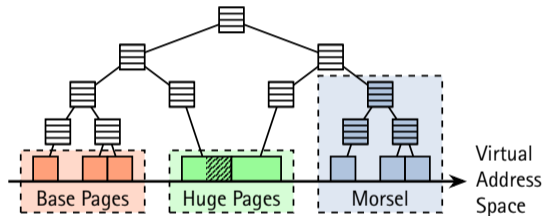


■ **Problem:** Conventional memory management

- Many individual entities (pages)
- Huge pages induce fragmentation
- No persistent state

■ **Solution:** Morsels

- Indivisible unit (page-table subtree)
- Ready for persistent memory
- Efficiently sharable



■ **Problem:** Conventional memory management

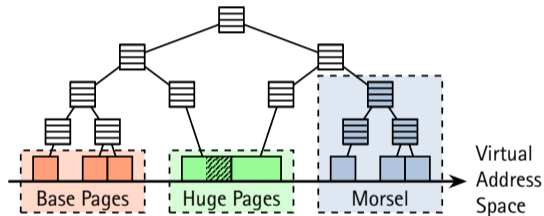
- Many individual entities (pages)
- Huge pages induce fragmentation
- No persistent state

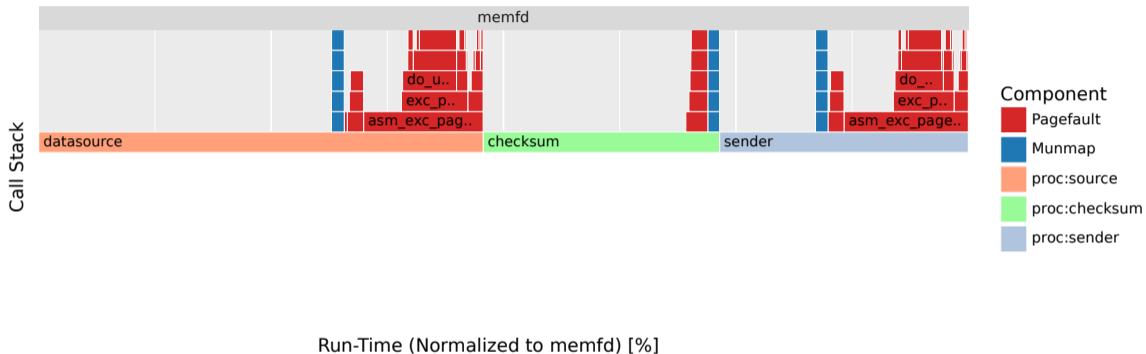
■ **Solution:** Morsels

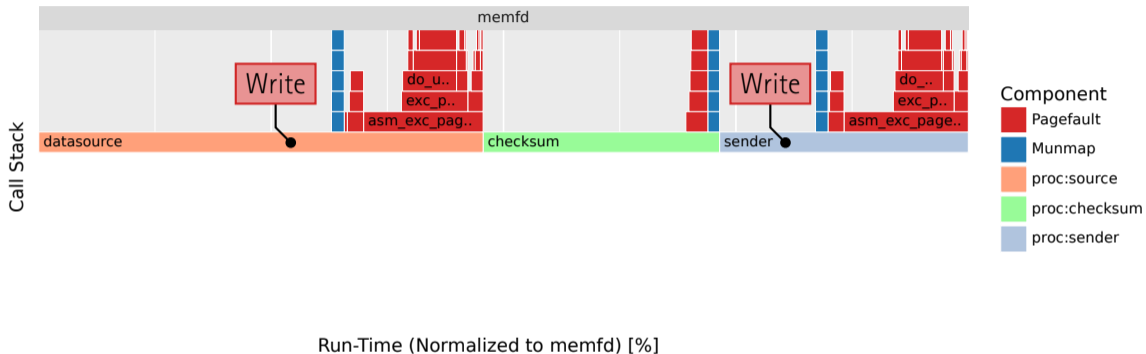
- Indivisible unit (page-table subtree)
- Ready for persistent memory
- Efficiently sharable

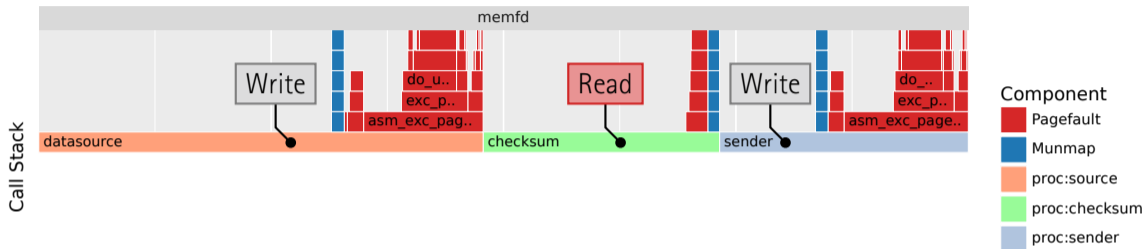
■ **Results**

- Constant mapping time (x1000 faster for 7 GiB)
- 39% reduced RTT for shown pipeline

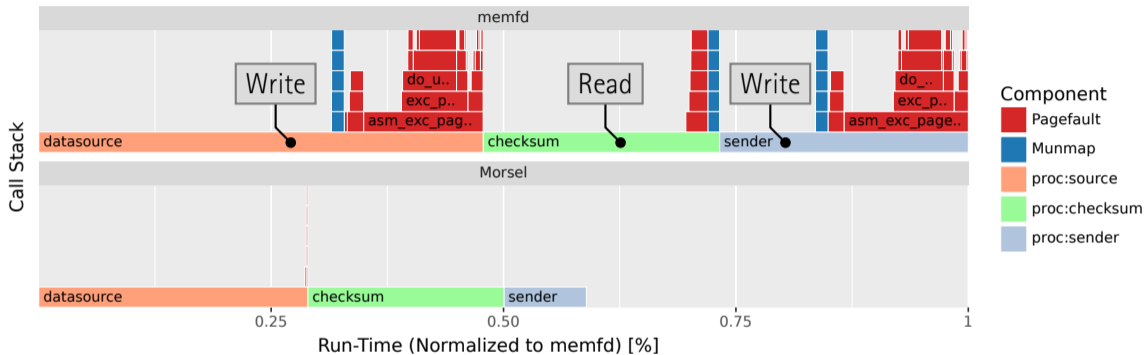


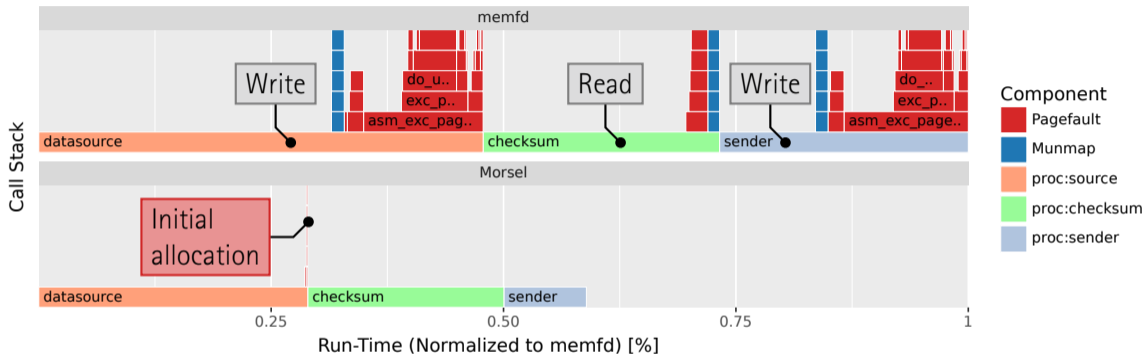


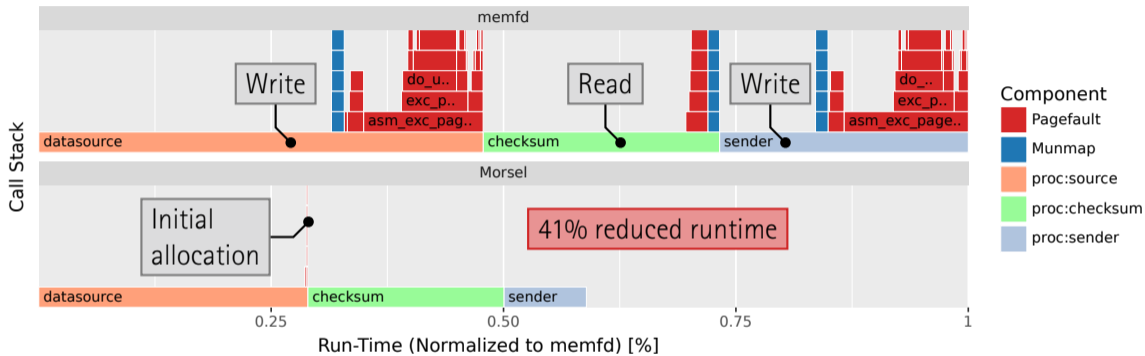


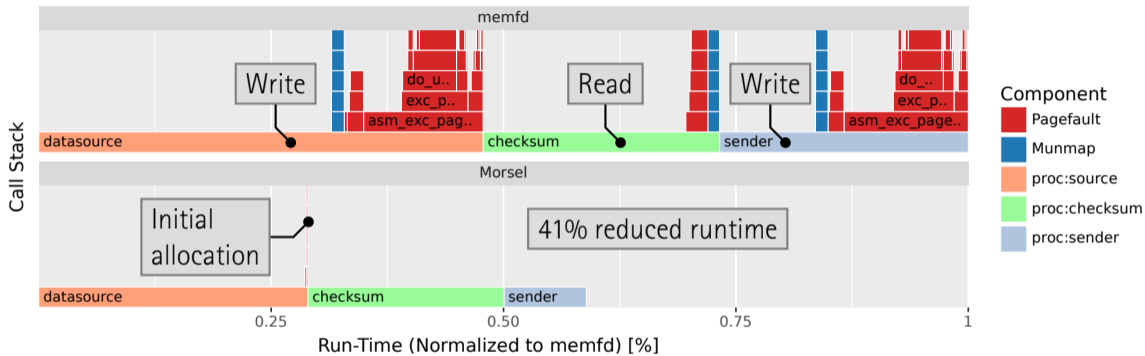


Run-Time (Normalized to memfd) [%]









Conclusion: Morsels enable efficient handling of virtual memory!