# The New Costs of Physical Memory Fragmentation

## DIMES '24

**Alexander Halbuer[1], Illia Ostapyshyn[1], Lukas Steiner[2], Lars Wrenger[1], Matthias Jung[3], Christian Dietrich[4], Daniel Lohmann[1]**

[1]Leibniz Universität Hannover
[2]Rheinland-Pfälzische Technische Universität Kaiserslautern-Landau
[3]Universität Würzburg and Fraunhofer IESE
[4]Technische Universität Braunschweig

November 03, 2024

- Motivation
  - One third of TCO and power consumption
  - Average utilization is at 75%

- **Motivation**
  - One third of TCO and power consumption
  - Average utilization is at 75%

- **Beliefs**
  - Only used memory is good memory (cache 'em all)
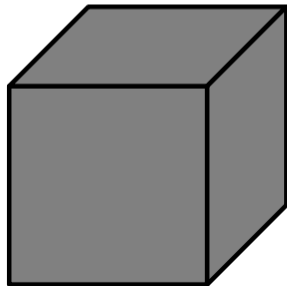  - All frames are equal (no external fragmentation)

- Motivation
  - One third of TCO and power consumption
  - Average utilization is at 75%

- Beliefs
  - Only used memory is good memory (cache 'em all)
  - All frames are equal (no external fragmentation)

- This is no longer the reality!
  - Huge/giant pages
  - Redistribute between VMs/via CXL
  - or power down unused memory

- Motivation
  - One third of TCO and power consumption
  - Average utilization is at 75%

- Beliefs
  - Only used memory is good memory (cache 'em all)
  - All frames are equal (no external fragmentation)

- This is no longer the reality!
  - Huge/giant pages
  - Redistribute between VMs/via CXL
  - or power down unused memory
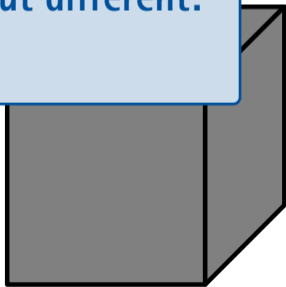
**4KiB**
**What everything is about**

# Physical Memory

- Motivation
  - One third of TCO and power consumption
  - Average utilization is at 75%

- Beliefs
  - Only used memory is good memory (cache 'em all)
  - All frames are equal (no external fragmentation)

- This is no longer the reality!
  - Huge/giant pages
  - Redistribute between VMs/via CXL
  - or power down unused memory

**4KiB**
What everything is about

**2MiB** / 1GiB
What we have to deal with

- Motivation
  - One third of TCO and power consumption
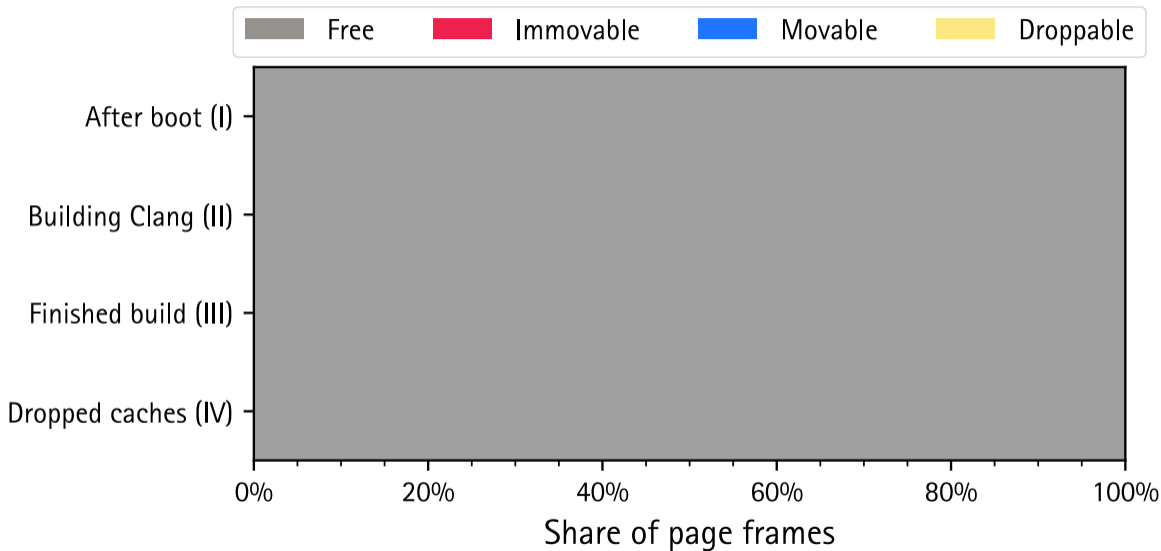  - Average utilization is at 75%

- B
  -
  -

- This is no longer the reality!
  - Huge/giant pages
  - Redistribute between VMs/via CXL
  - or power down unused memory
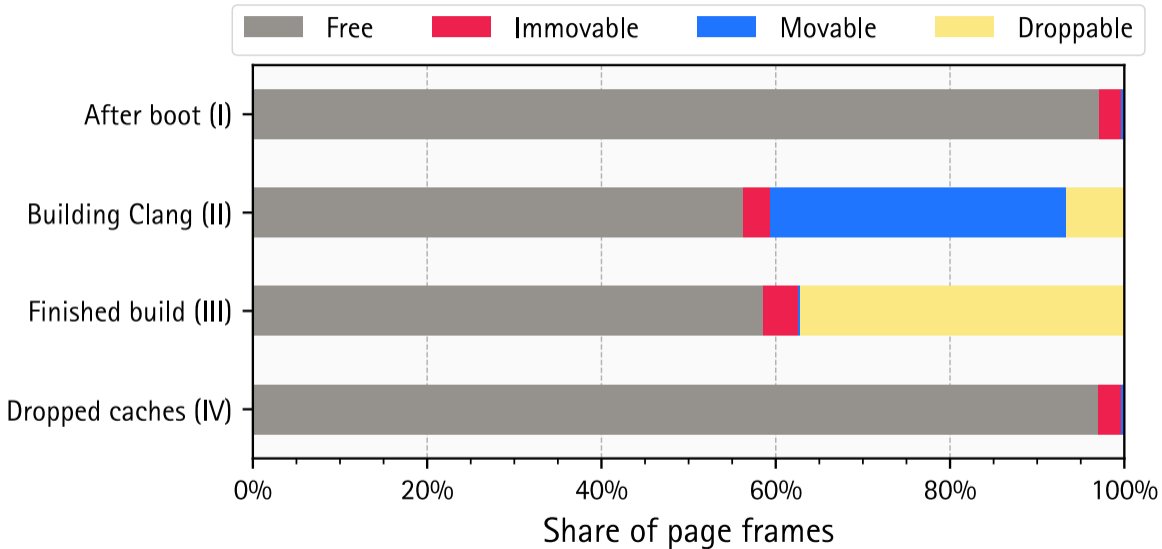
**4KiB**
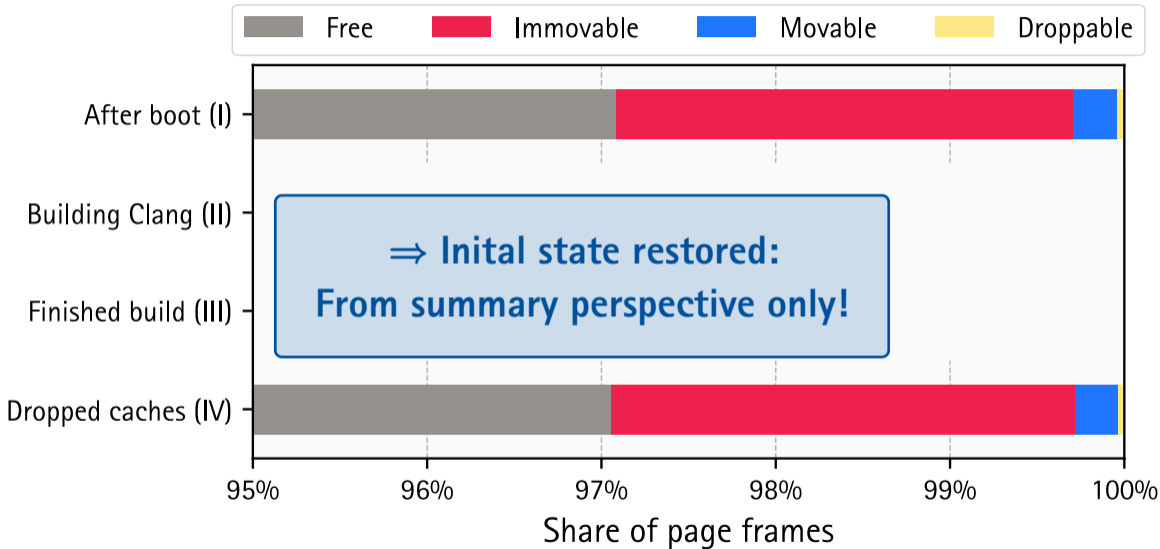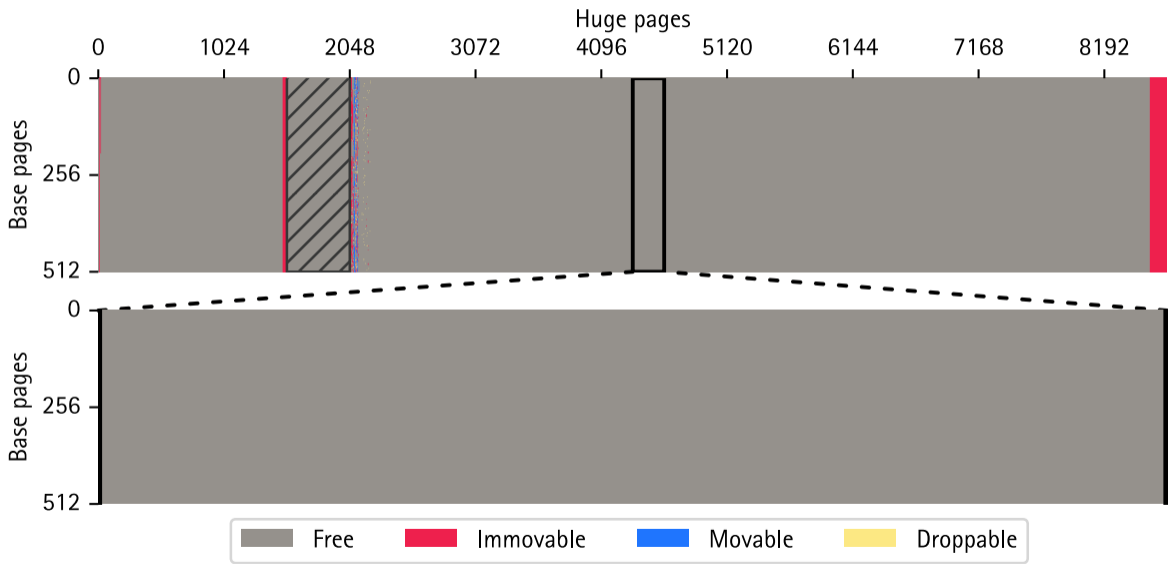What everything is about

**2MiB** / **1GiB**
What we have

⇒ **The new fragmentation problem is real but different:**
**We must think bigger!**

Case Study: Building Clang, 12 CPUs, 16 GiB

Legend: Free, Immovable, Movable, Droppable

Categories (top to bottom):
- After boot (I)
- Building Clang (II)
- Finished build (III)
- Dropped caches (IV)

X-axis: Share of page frames (0% to 100%)

# Case Study: Building Clang, 12 CPUs, 16 GiB

Case Study: Building Clang, 12 CPUs, 16 GiB

Legend: Free (gray), Immovable (red), Movable (blue), Droppable (yellow)

Categories: After boot (I), Building Clang (II), Finished build (III), Dropped caches (IV)

X-axis: Share of page frames (95%, 96%, 97%, 98%, 99%, 100%)

⇒ **Inital state restored:
From summary perspective only!**

(I) After boot — (IV) Dropped caches

Free blocks [%] (more is better) vs Block size (4 KiB, 32 KiB, 256 KiB, 2 MiB, 16 MiB, 128 MiB, 1 GiB, 8 GiB)

Baseline

Δ=10 blocks

Cost function: $\#\text{TOUCHes} = \#\text{DROPs} + 2 \times \#\text{MOVEs}$

Cost function: $\#\text{TOUCHes} = \#\text{DROPs} + 2 \times \#\text{MOVEs}$

Cost function: $\#\text{TOUCHes} = \#\text{DROPs} + 2 \times \#\text{MOVEs}$



$\Rightarrow$ **Immovable pages lead to big lost potential!**
- **Compaction amortizes in < 1min (power off)**
- **9 more GiB unused but not contiguous**

Legend:
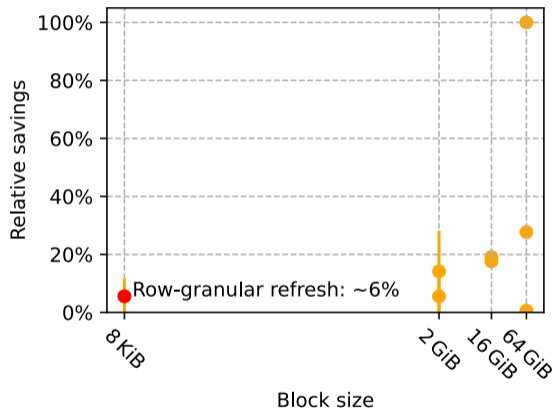- OPT: move everything (incl. immovable)
- Max compaction
- 1% frames touched
- Baseline

Left axis: Free b (more) — values 20, 40
X-axis (Block size): 4 KiB, 32 KiB, 256 KiB, 2 MiB, 16 MiB, 128 MiB, 1 GiB, 8 GiB

Right chart annotations: $\Delta$=3 blocks, $\Delta$=6 blocks, $\Delta$=4 blocks

# DRAM Power Saving



**Example system**

512 GiB DDR5

8 channels

4 ranks per channels

8×2 GiB chips per rank

32 banks per rank

65 536 rows per bank

# DRAM Power Saving



**Example system**

512 GiB DDR5

8 channels

4 ranks per channels

8×2 GiB chips per rank

32 banks per rank

65 536 rows per bank

# DRAM Power Saving



**Example system**
512 GiB DDR5
8 channels
4 ranks per channels
8×2 GiB chips per rank
32 banks per rank
65 536 rows per bank

# DRAM Power Saving
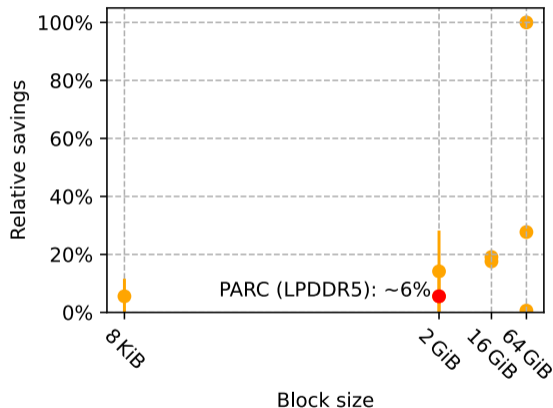
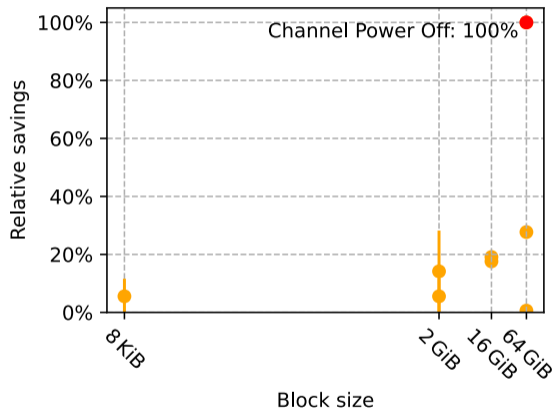

**Example system**

512 GiB DDR5

8 channels

4 ranks per channels

8×2 GiB chips per rank

32 banks per rank

65 536 rows per bank
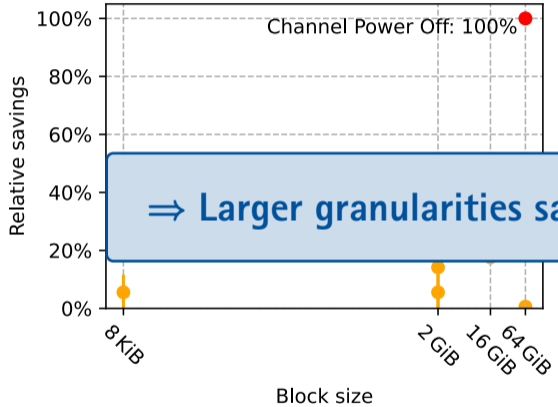
**Example system**

512 GiB DDR5

8 channels

4 ranks per channels

... per rank

... ank

65 536 rows per bank

⇒ **Larger granularities save more power**

- **Linux is bad at managing physical memory**

  - No fragmentation avoidance above (2 MiB) huge page size

  - Page cache unconditionally fills all available memory

  - Immovable pages make compaction impossible

- **Using memory comes at a cost**

  - Huge/giant page availability

  - Redistribute or turn off unused memory

  - Dynamic cloud pricing models

https://sra.uni-hannover.de/p/dram-dimes24