# WANT MORE UNIKERNELS? INFLATE THEM!

**GAULTHIER GAIN**[1], CYRIL SOLDANI[1], FELIPE HUICI[2], PROF. LAURENT MATHY[1]

[1]University of Liège, [2]Unikraft GmbH

LIÈGE université

Unikraft

# DILEMMA

**Virtual Machines (VMs)**

✔ Strong isolation

✘ Heavyweight
   – Degrade performance

**Containers**

✘ Poor isolation
   – A lot of exploits

✔ Lightweight
   – Share underlying kernel

# DILEMMA

**Virtual Machines (VMs)**
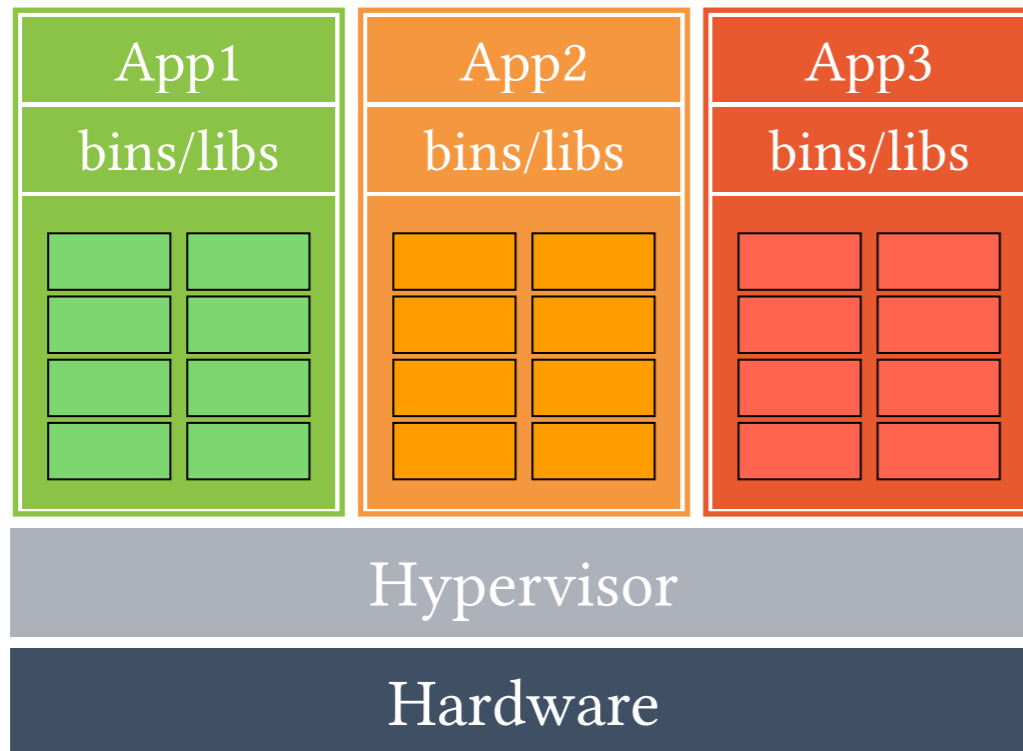
✔️ Strong isolation

❌ Heavyweight
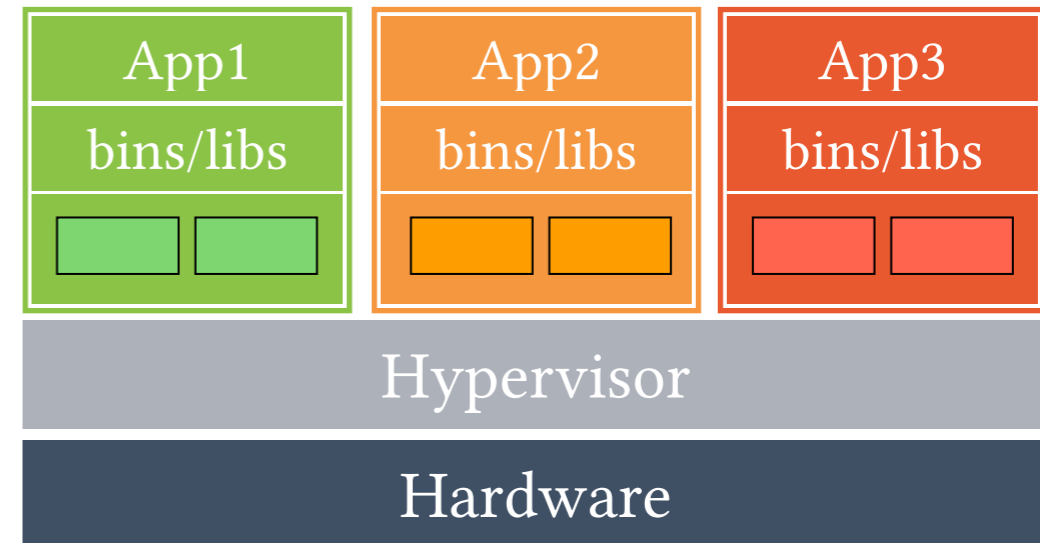– Degrade performance

**Containers**

❌ Poor isolation
– A lot of exploits

✔️ Lightweight
– Share underlying kernel

## Solution → **Unikernels**

# UNIKERNELS
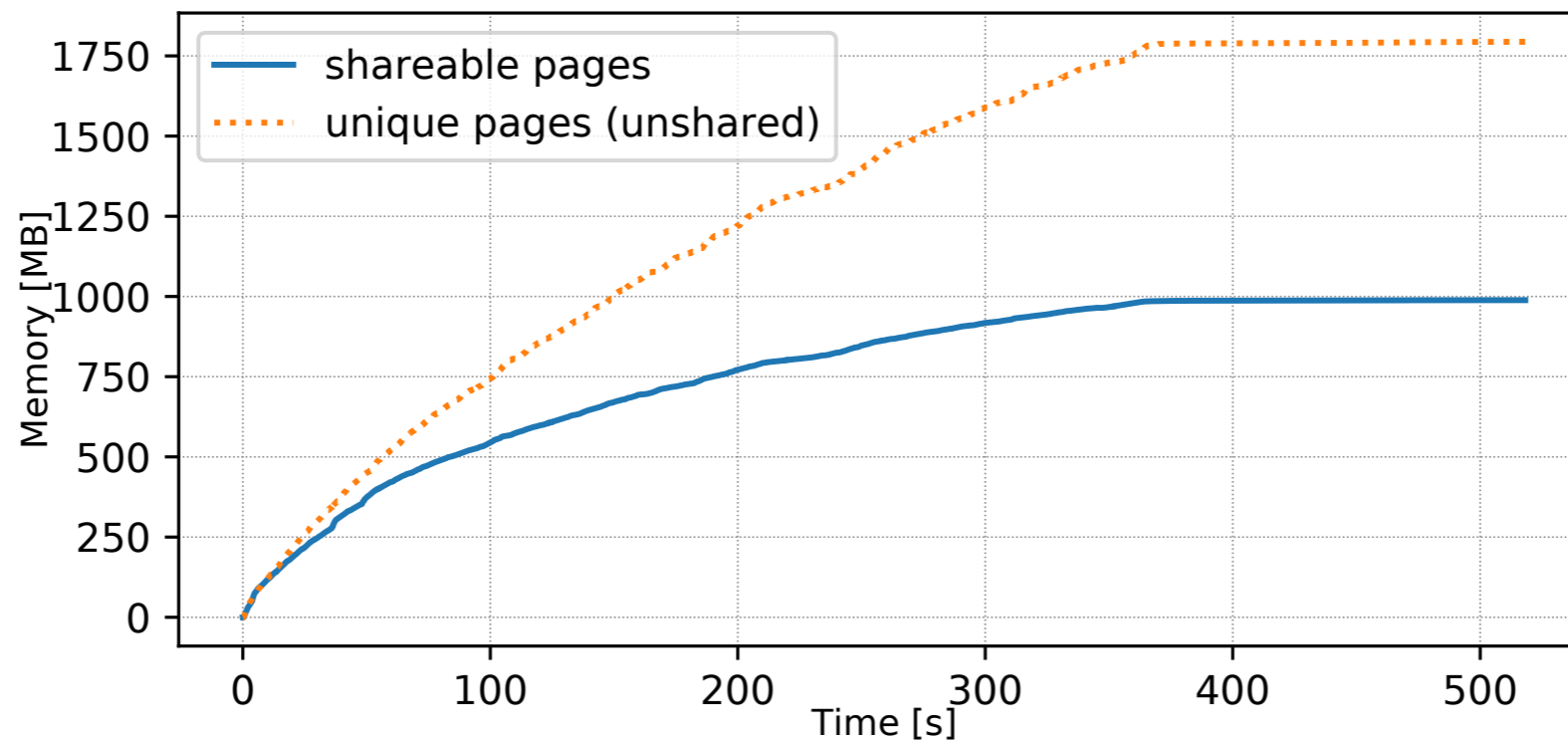


Virtual Machines (VMs)                    **Unikernels**

Unikernels are purpose-built:

▸ Thin kernel layer (only the necessary features that the application needs).

▸ Essential functions are placed into micro-libs (µlibs) with well-defined behaviour.

# UNIKERNELS GAINS

▸ Fast instantiation, destruction and migration times:

  ▸ Hundred of milliseconds.

▸ Small per-instance memory footprint:

  ▸ Few MBs or even KBs.

▸ High performance:

  ▸ 10-40 Gbps throughput.

▸ Reduced attack surface:

  ▸ Less components.

▸ High density:

  ▸ Thousand of instances on a single host $\rightarrow$ Can we do better ?
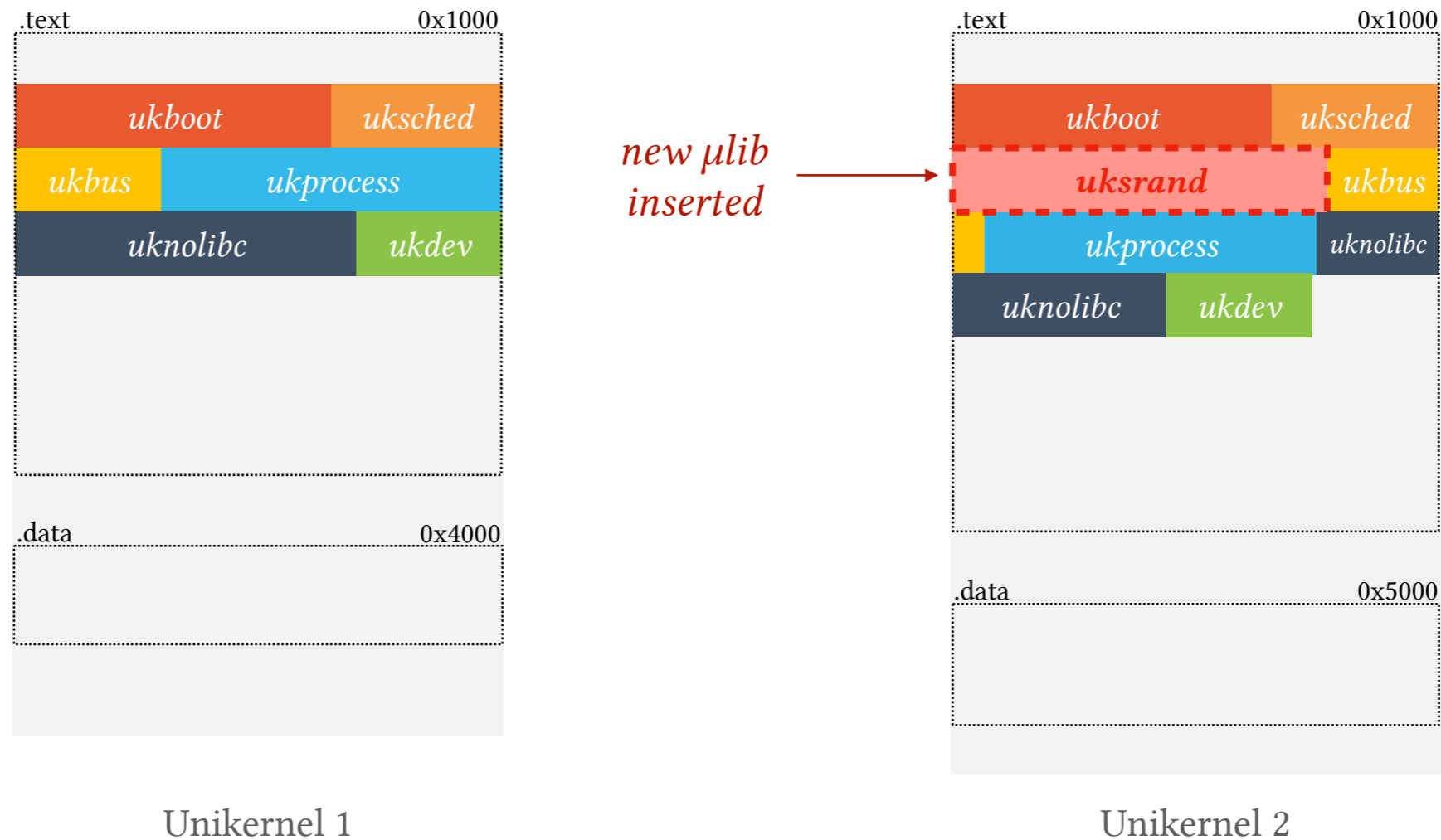
# RUNNING A LARGE NUMBER OF UNIKERNELS



Evolution of unshared and shareable (i.e. having at least one copy) pages when running 1000 different FaaS unikernels (with ASLR) on a single physical server.
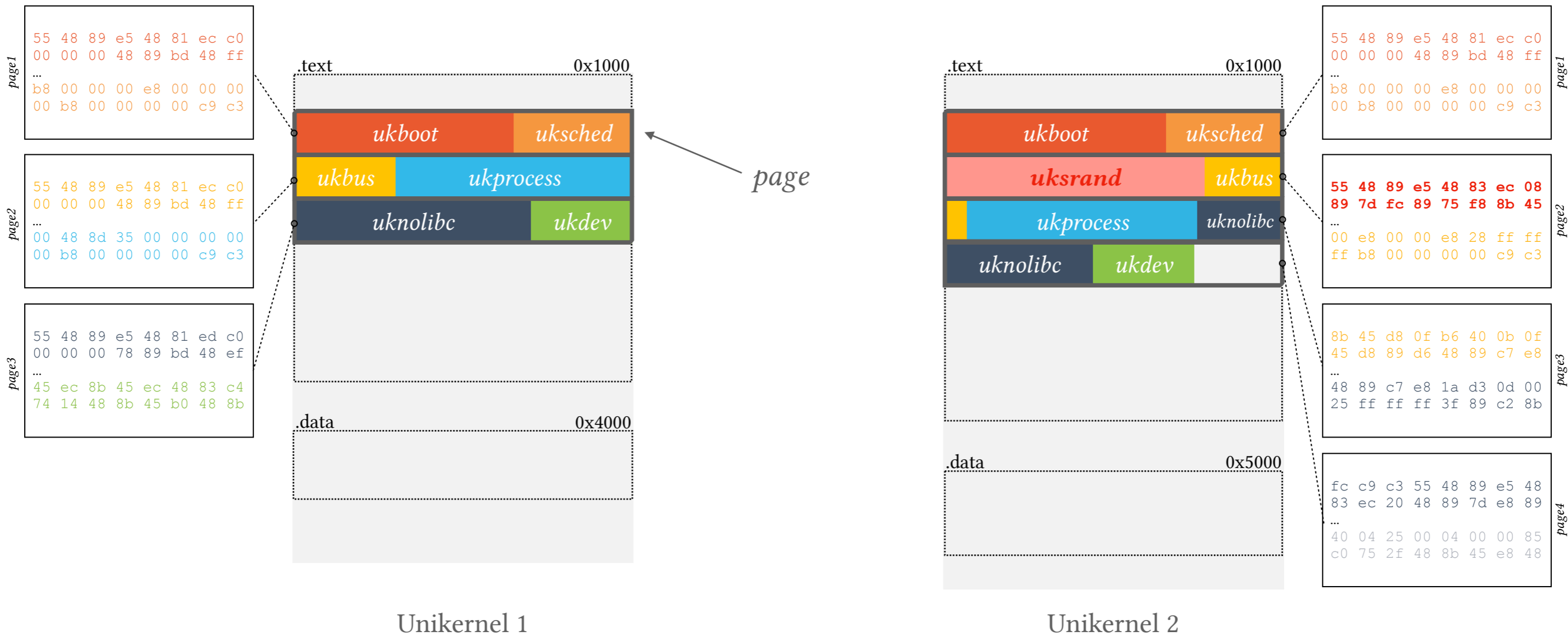
▸ We investigated by running a large number of unikernels on a same physical server and we relied on a memory deduplication scanner (UKSM*).

▸ Unique pages are much more frequent than shared pages.

▸ Specialisation? Need further investigation to understand the reason.

*Ultra Kernel Samepage Merging

# MEMORY DEDUPLICATION WITH UNIKERNELS: OVERVIEW



Unikernel 1                    Unikernel 2

▸ Having several instances will result into different µlibs configurations.

▸ The underlying build system does not have a global overview of the µlibs: Each unikernel is built in an individual way.

▸ All µlibs are compacted: resulting unikernel consumes as little memory and disk space as possible.

Unikernel 1

Unikernel 2

If a new μlib '*uksrand*' is inserted between other μlibs:

▸ μlibs' code will be split across different pages.

▸ It reduces memory sharing since pages are different.
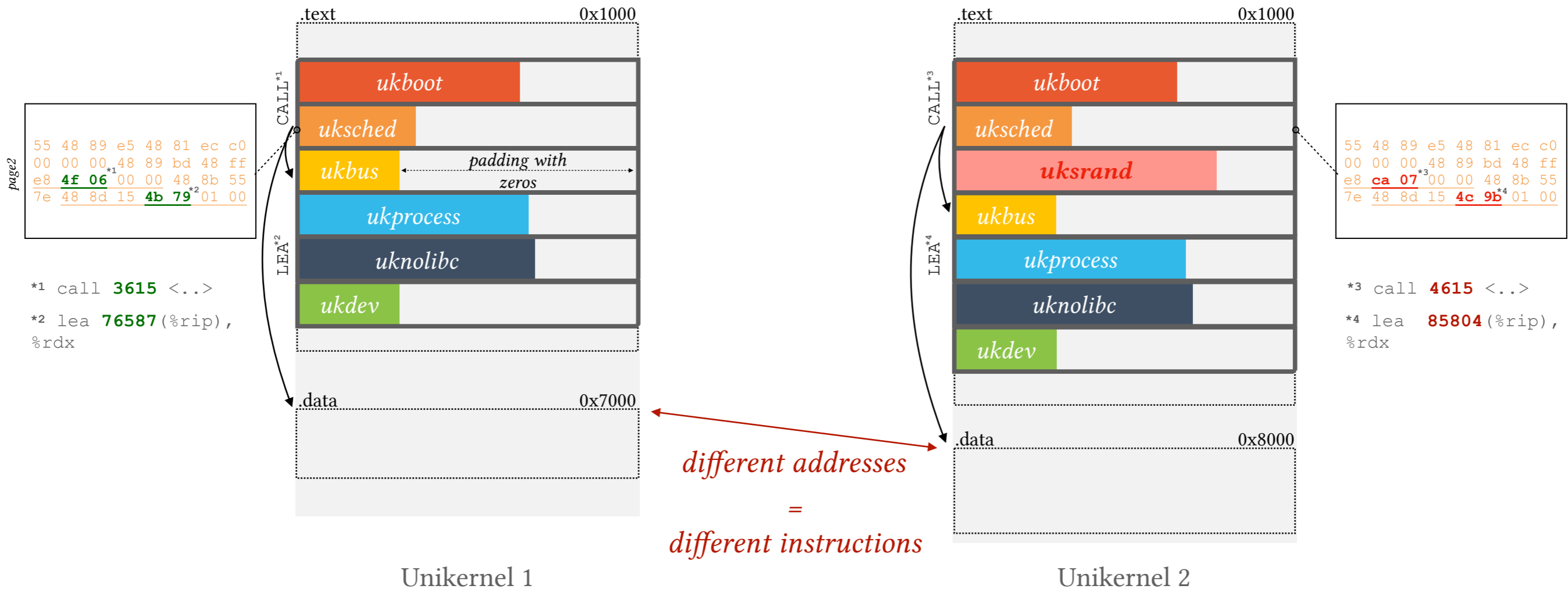
# MEMORY DEDUPLICATION WITH UNIKERNELS: A FIRST SOLUTION?



Unikernel 1

Unikernel 2

To circumvent this issue:

▸ Align each μlib to a page boundary address.

▸ Pad the μlib code with zeros to fill a complete page.

→ Is it enough?

9

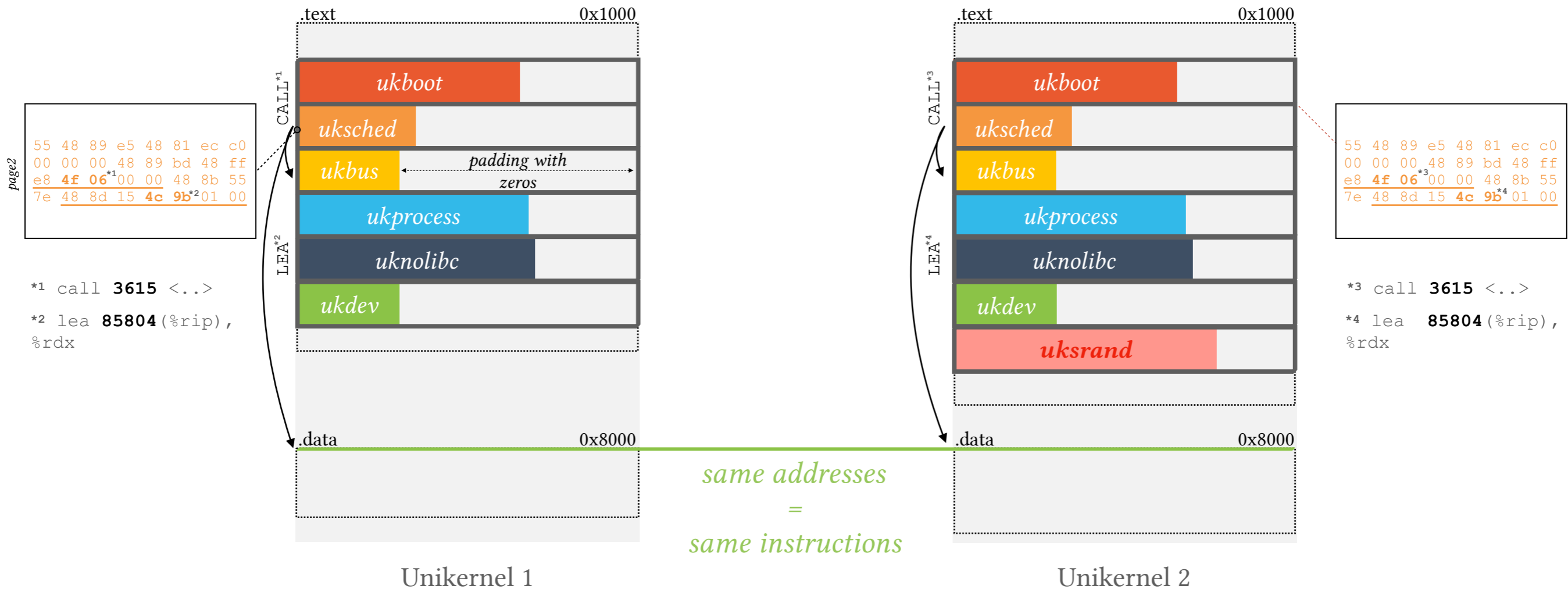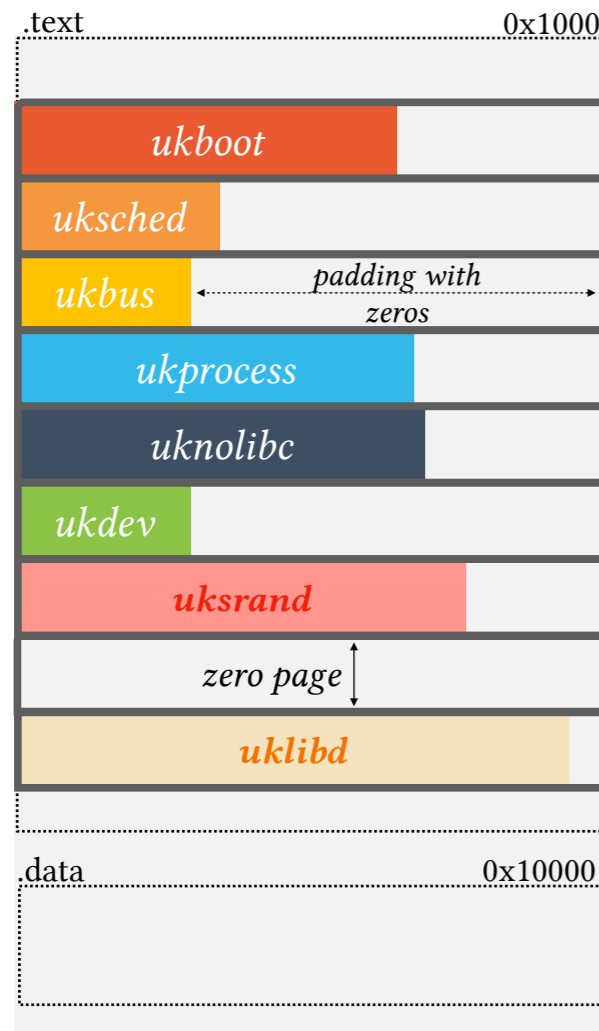Some instructions use different addresses in the *.text* section:

▸ Related to other sections (e.g., *.data*, *.rodata*): `MOV` and `LEA`.

▸ Related to another part of the *.text* section: `CALL`.

# MEMORY DEDUPLICATION WITH UNIKERNELS:  A WORKING SOLUTION
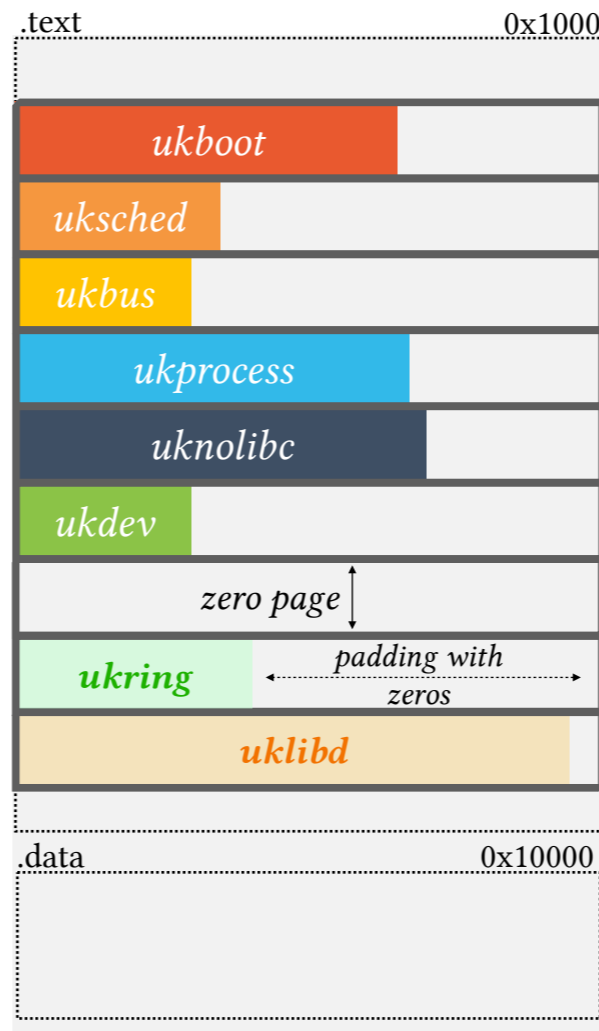


Unikernel 1                    Unikernel 2

same addresses
=
same instructions

1. Placing μlibs at page boundary addresses.

2. Keep a same μlibs order.

3. Align sections (e.g., *.data, .rodata, ...*) at same addresses.

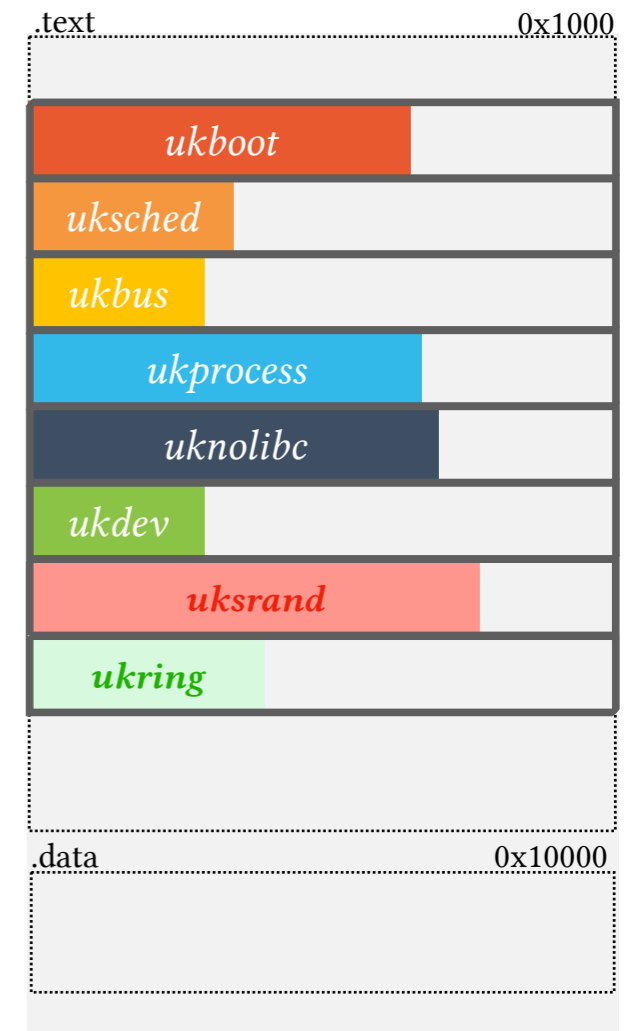# MEMORY DEDUPLICATION WITH UNIKERNELS



Unikernel 1      Unikernel 2      Unikernel 3

If there are more than two instances with different μlibs subsets:
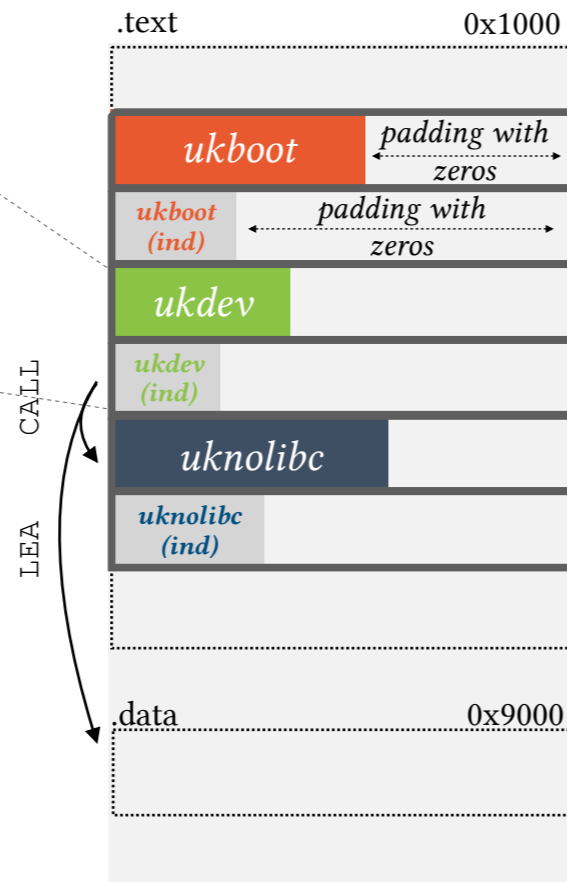
▸ It is necessary to align them to specific addresses.

▸ This leads to 'gaps' of zero pages in the memory virtual space.

# TOWARDS ASLR SUPPORT



1. Relative jump to indirection table;
2. Execute the problematic instruction (indirection);
3. Jump back to the microlib code.

```
ukdev (0x2100):
    jmp eip+1100(0x3200)
    …
ukdev(ind) (0x3200):
    call 0x4200
    jmp 0x2105
```

```
ukdev (0x4100):
    jmp eip+1100(0x5200)
    …
ukdev(ind) (0x5200):
    call 0x1200
    jmp 0x4105
```

Unikernel1

Unikernel2

▸ Using fixed absolute addresses leads to security issues (no ASLR).

▸ Create an indirection table per μlib which contains problematic instructions (using addresses from other sections/μlibs). Such instructions are replaced by relative jump to their new position.

# SPACER HIGH-LEVEL ARCHITECTURE

*1 These steps are only performed in Spacer (ASLR)

workspace

| | | | |
|---|---|---|---|
| unikernel 1 | unikernel 2 | … | unikernel N |

① Analyse unikernel

④ Binary rewriting*1

*Spacer*

Unikraft

② Generate linker.lds
(libs are shuffled)*1

③ Linking

▸ From our methodology, we derive Spacer, a tool aims to have a global knowledge of all the µlibs used by all unikernels on the same workspace.

▸ Spacer performs a new linking by associating µlibs with absolute addresses according to a map (by rewriting the linker script).

▸ For Spacer (ASLR), µlibs are shuffled during the linker file generation. Furthermore, there is one extra step of binary rewriting (move problematic instructions to indirection tables).

LIÈGE université

# EVALUATION: METHODOLOGY

► We compared Spacer with DCE (Dead Code Elimination) and Default configuration.

  ► 10 applications ported as unikernels.

  ► 1000 FaaS unikernels.

  ► On several dimensions: memory consumption, file size and performance.

*Please refer to the paper for the full evaluation.*

# EVALUATION (1)

Memory consumption:

- ▸ Without memory deduplication, Spacer and Spacer (ASLR) consume significantly more memory (zero pages and indirection tables).

- ▸ With memory deduplication, the benefits of alignment increases as we run more applications. Spacer and Spacer (ASLR) consume less memory than default and DCE.

- ▸ Up to a 3x gain compared to DCE.

*Heap-intensive applications*:

- ▸ The gain is less noticeable (e.g., in-memory databases).

- ▸ If there are thousands of applications, Spacer still allows to reduce the memory consumed (code and read-only data are shared).

- ▸ But if they are only some instances: do not apply Spacer on it.

# EVALUATION (2)

Elf Size:

- Spacer and Spacer (ASLR) have a slight impact on file size:

    - The inflation of the header string table (ELF section).

    - Indirection tables (problematic instructions).

- Spacer ELF files do not have inflation due to zeros, it is only in memory.

Performance:

- Total execution time of short-lived and long-lived unikernels.

- UKSM has a slight impact on scanning and merging pages.

- Spacer performance degradation is minimal: having zero pages and indirection tables introduces a slight overhead (± 5-7% compared to DCE).

LIÈGE université

# CONCLUSION & FUTURE WORK

▸ Unikernels are small and have impressive performance, but they show few opportunities for VM page sharing (specialisation).

▸ We brought a new methodology that rearranges and inflates unikernels by using μlibs alignment.

  ▸ Aligning μlibs may lead up to a big reduction in memory consumption, even when compared to unikernels built with DCE (Dead Code Elimination).

  ▸ Furthermore, the alignment does not introduce significant overhead in terms of ELF size, nor does it impairs application performance.

Future work:

  ▸ *Loader*: A loader that performs deduplication at load time could make μlib pages point directly to the corresponding frames when loading the kernel image into main memory.

# THANK YOU!

# QUESTIONS?

# SPACER & UNIKRAFT

Spacer: https://people.montefiore.uliege.be/gain/spacer

Unikraft: https://unikraft.org/

Contact by email: gaulthier.gain@uliege.be

License: *3-Clause BSD License*

Gaulthier Gain, Cyril Soldani, Felipe Huici, Prof. Laurent Mathy