

# 内存管理的层次结构

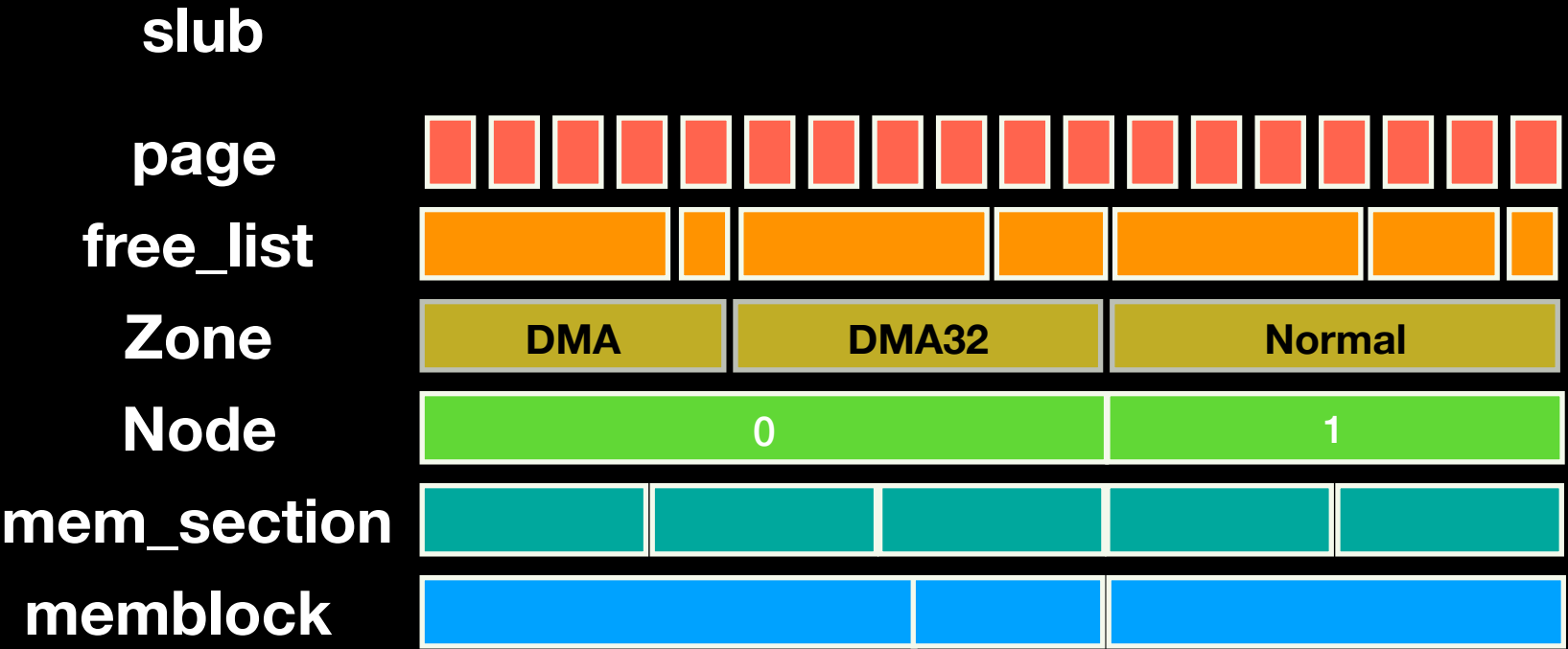
Wei Yang

<richard.weiyang@gmail.com>

# 议程

- 现有的层次结构
- 分层的原因

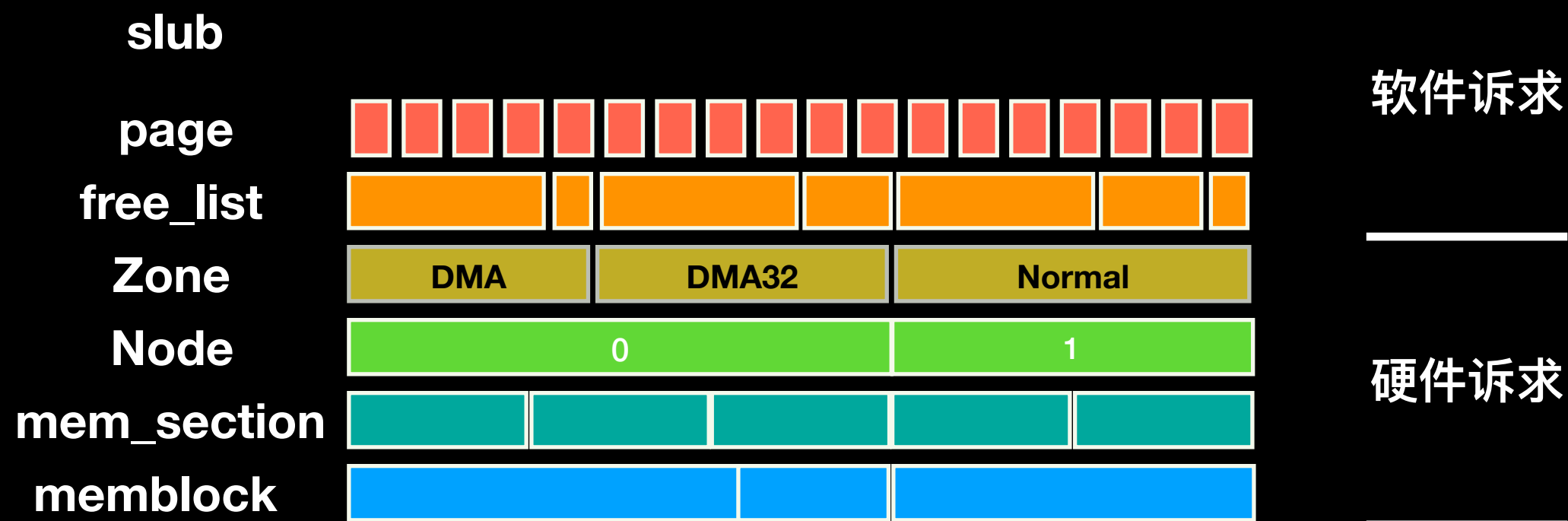
# 现有的层次结构



内核是软件和硬件的结合  
是两者妥协的产物

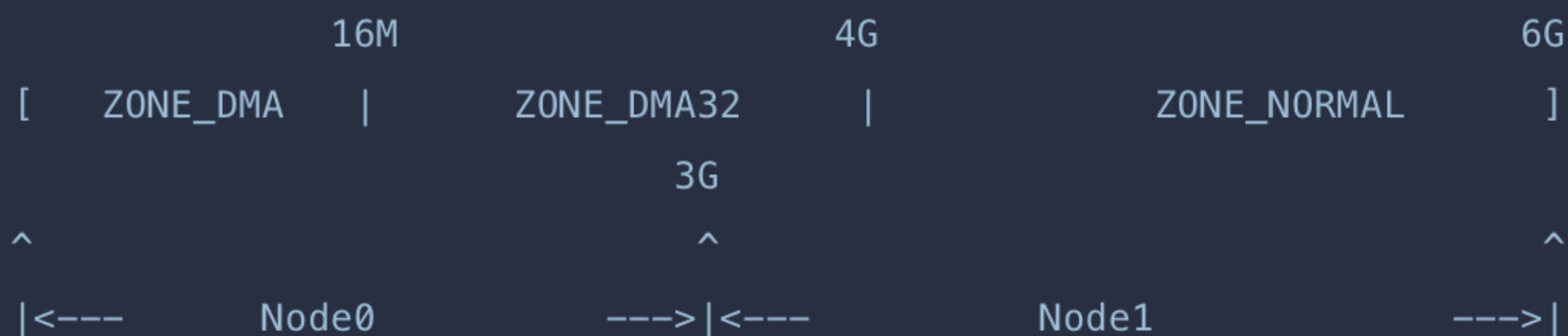
# 分层的原因

- 描述硬件属性： Node, Zone, Page
- 特性支持： hotplug
- 提升扩展性： zone->free\_list, slub, pcp
- 节省空间： mem\_section



# Node 和 Zone的含义

Memory



# 热插拔支持





# 扩展性考量 1

```
node_data[0]
+-----+
|node_id          <---+ |
|  (int)           | |
+-----+
|node_zones[MAX_NR_ZONES] | | [ZONE_DMA]
|  (struct zone)         | | +-----+
|  +-----+             | | |0          |
|  |                   | | |16M         |
|  |zone_pgdat         ----+ | +-----+
|  |                   | |
|  |                   | [ZONE_DMA32]
|  |                   | +-----+
|  |                   | |16M         |
|  |                   | |3G          |
|  |                   | +-----+
|  |                   |
|  |                   | [ZONE_NORMAL]
|  |                   | +-----+
|  |                   | |empty       |
|  |                   | |
+-----+ +-----+
```

# 扩展性考量 2

```
struct zone
```

```
+-----+
|pageset|
|  (struct per_cpu_pageset *)|
|  cpu0          cpu1          ...    cpuN|
|  +-----+ +-----+ ... +-----+
|  |pcp|      |pcp|      |pcp|
|  | (struct per_cpu_pages)| | (struct per_cpu_pages)| | (struct per_cpu_pages)|
|  | +-----+ | +-----+ | +-----+
|  | |count|   | |count|   | |count|
|  | |high|    | |high|    | |high|
|  | |batch|   | |batch|   | |batch|
|  | |        | | |        | | |
|  | |lists[MIGRATE_PCPTYPES]| | |lists[MIGRATE_PCPTYPES]| | |lists[MIGRATE_PCPTYPES]|
+---+---+-----+---+---+-----+---+---+-----+
```

# 节省空间 page的存放

```
mem_section[NR_SECTION_ROOTS][SECTIONS_PER_ROOT]
```

```
= [DIV_ROUND_UP(NR_MEM_SECTIONS, SECTIONS_PER_ROOT)][SECTIONS_PER_ROOT]
```

	[0]	[1]				[SECTIONS_PER_ROOT - 1]
	+-----+-----+				+-----+-----+	
[0]			...			
	+-----+-----+				+-----+-----+	
	+-----+-----+				+-----+-----+	
[1]			...			
	+-----+-----+				+-----+-----+	
	+-----+-----+				+-----+-----+	
[2]			...			
	+-----+-----+				+-----+-----+	