# NUMA -1- (ARM64 initialization)

📅 2019-04-19 (http://jake.dothome.co.kr/numa-1/)    👤 Moon Young-il
(http://jake.dothome.co.kr/author/admin/)    📂 Linux Kernel (http://jake.dothome.co.kr/category/linux/)

<kernel v5.10>

## NUMA -1- (ARM64 initialization)

### arm64_numa_init()

arch/arm64/mm/numa.c

```
1  /**
2   * arm64_numa_init - Initialize NUMA
3   *
4   * Try each configured NUMA initialization method until one succeeds.  T
    he
5   * last fallback is dummy single node config encomapssing whole memory.
6   */
```

```
01  void __init arm64_numa_init(void)
02  {
03          if (!numa_off) {
04                  if (!acpi_disabled && !numa_init(arm64_acpi_numa_init))
05                          return;
06                  if (acpi_disabled && !numa_init(of_numa_init))
07                          return;
08          }
09
10          numa_init(dummy_numa_init);
11  }
```

On an ARM64 system, attempt to configure a NUMA system through the ACPI table or device tree. If there is no such configuration, it performs a dummy numa initialization with a single node.

### numa_init()

arch/arm64/mm/numa.c

```
01  static int __init numa_init(int (*init_func)(void))
02  {
03          int ret;
04
05          nodes_clear(numa_nodes_parsed);
06          nodes_clear(node_possible_map);
07          nodes_clear(node_online_map);
08
09          ret = numa_alloc_distance();
10          if (ret < 0)
11                  return ret;
12
13          ret = init_func();
14          if (ret < 0)
15                  goto out_free_distance;
16
```

```
17          if (nodes_empty(numa_nodes_parsed)) {
18                  pr_info("No NUMA configuration found\n");
19                  ret = -EINVAL;
20                  goto out_free_distance;
21          }
22
23          ret = numa_register_nodes();
24          if (ret < 0)
25                  goto out_free_distance;
26
27          setup_node_to_cpumask_map();
28
29          return 0;
30  out_free_distance:
31          numa_free_distance();
32          return ret;
33  }
```

Initialize the NUMA system.

- In line 5~7 of the code, initialize all bitmaps that manage the NUMA state.
- In code lines 9~11, allocate as many numa_disatance[] arrays as you need and specify the initial values.
- Run the function @init_func received as an argument in line 13~15 of the code to set it.
- If there is no node information in line 17~21, it prints the message "No NUMA configuration found" and exits.
- In lines 23~25 of the code, assign the node data (struct pglist_data) and set the default node information.
- In line 27 of code, set the node-> cpu map.

As shown below, if the NUMA setting cannot be found through ACPI or the device tree, you will see the message "No NUMA configuration found" output. And it has been fallback, and through the dymmy NUMA initialization routine it says "Faking a node at..." You can also see the message being printed.

```
01  [    0.000000] NUMA: No NUMA configuration found
02  [    0.000000] NUMA: Faking a node at [mem 0x0000000040000000-0x00000000
    7fffffff]
03  [    0.000000] NUMA: NODE_DATA [mem 0x7ddf4840-0x7ddf5fff]
04  [    0.000000] Zone ranges:
05  [    0.000000]   DMA32    [mem 0x0000000040000000-0x000000007fffffff]
06  [    0.000000]   Normal   empty
07  [    0.000000] Movable zone start for each node
08  [    0.000000] Early memory node ranges
09  [    0.000000]   node   0: [mem 0x0000000040000000-0x000000007fffffff]
10  [    0.000000] Initmem setup node 0 [mem 0x0000000040000000-0x000000007f
    ffffff]
```

## numa_alloc_distance()

arch/arm64/mm/numa.c

```
01  static int __init numa_alloc_distance(void)
02  {
03          size_t size;
04          u64 phys;
05          int i, j;
06
```

```
07          size = nr_node_ids * nr_node_ids * sizeof(numa_distance[0]);
08          phys = memblock_find_in_range(0, PFN_PHYS(max_pfn),
09                                        size, PAGE_SIZE);
10          if (WARN_ON(!phys))
11                  return -ENOMEM;
12
13          memblock_reserve(phys, size);
14
15          numa_distance = __va(phys);
16          numa_distance_cnt = nr_node_ids;
17
18          /* fill with the default distances */
19          for (i = 0; i < numa_distance_cnt; i++)
20                  for (j = 0; j < numa_distance_cnt; j++)
21                          numa_distance[i * numa_distance_cnt + j] = i ==
j ?
22                                  LOCAL_DISTANCE : REMOTE_DISTANCE;
23
24          pr_debug("Initialized distance table, cnt=%d\n", numa_distance_c
nt);
25
26          return 0;
27  }
```

Number of nodes * Create and allocate an array of numa_distance[] for the number of nodes.

- In line 7 of code, use the number of nodes * the number of nodes as the size to be assigned.
- In line 8~11 of the code, we find out the space equal to the size of the empty space of the memblock on a page-by-page basis.
- In line 13 of code, reserve the size area to the memblock.
- Convert the physical address assigned in line 15~16 to a virtual address, assign it to a numa_distance, and specify the number of nodes.
- In lines 19~22 of code, numa_distance[] logically expresses a double array. The default numa_distance[] value is LOCAL_DISTANCE(10) only if it means the same node, and REMOTE_DISTANCE(20) is otherwise specified.
  - e.g. if nr_node_ids=2, look at numa_distance[].
    - numa_distance[0] = 10
    - numa_distance[1] = 20
    - numa_distance[2] = 20
    - numa_distance[3] = 10

## numa_register_nodes()

arch/arm64/mm/numa.c

```
01  static int __init numa_register_nodes(void)
02  {
03          int nid;
04          struct memblock_region *mblk;
05
06          /* Check that valid nid is set to memblks */
07          for_each_mem_region(mblk) {
08                  int mblk_nid = memblock_get_region_node(mblk);
09
10                  if (mblk_nid == NUMA_NO_NODE || mblk_nid >= MAX_NUMNODE
S) {
11                          pr_warn("Warning: invalid memblk node %d [mem %#
010Lx-%#010Lx]\n",
```

```
12                              mblk_nid, mblk->base,
13                              mblk->base + mblk->size - 1);
14                      return -EINVAL;
15              }
16      }
17
18      /* Finally register nodes. */
19      for_each_node_mask(nid, numa_nodes_parsed) {
20              unsigned long start_pfn, end_pfn;
21
22              get_pfn_range_for_nid(nid, &start_pfn, &end_pfn);
23              setup_node_data(nid, start_pfn, end_pfn);
24              node_set_online(nid);
25      }
26
27      /* Setup online nodes to actual nodes*/
28      node_possible_map = numa_nodes_parsed;
29
30      return 0;
31 }
```

Assign node data (struct pglist_data) and set the default node information.

- In code lines 7~16, if no node is specified in the memblock, it will print a warning message and return an error -EINVAL value.
- In line 19~25 of the code, traverse the nodes, assign node data (struct pglist_data), set the basic information of the node, and change the node status to online.
- In line 28 of code, specify the online nodes in the possible map.


## setup_node_to_cpumask_map()

arm64/mm/numa.c

```
1  /*
2   * Allocate node_to_cpumask_map based on number of available nodes
3   * Requires node_possible_map to be valid.
4   *
5   * Note: cpumask_of_node() is not valid until after this is done.
6   * (Use CONFIG_DEBUG_PER_CPU_MAPS to check this.)
7   */
```

```
01 static void __init setup_node_to_cpumask_map(void)
02 {
03      int node;
04
05      /* setup nr_node_ids if not done yet */
06      if (nr_node_ids == MAX_NUMNODES)
07              setup_nr_node_ids();
08
09      /* allocate and clear the mapping */
10      for (node = 0; node < nr_node_ids; node++) {
11              alloc_bootmem_cpumask_var(&node_to_cpumask_map[node]);
12              cpumask_clear(node_to_cpumask_map[node]);
13      }
14
15      /* cpumask_of_node() will now work */
16      pr_debug("Node to cpumask map for %u nodes\n", nr_node_ids);
17 }
```

Node - > sets the CPU map.

- In line 6~7 of the code, specify the number of possible nodes if the nr_node_ids has not yet been set.
- In line 10~13 of code, it iterates through the number of nodes, allocates and initializes the node_to_cpumask_map[node] bitmap.

## dummy_numa_init()

arch/arm64/mm/numa.c

```
1  /**
2   * dummy_numa_init - Fallback dummy NUMA init
3   *
4   * Used if there's no underlying NUMA architecture, NUMA initialization
5   * fails, or NUMA is disabled on the command line.
6   *
7   * Must online at least one node (node 0) and add memory blocks that cov
   er all
8   * allowed memory. It is unlikely that this function fails.
9   */
```

```
01  static int __init dummy_numa_init(void)
02  {
03          phys_addr_t start = memblock_start_of_DRAM();
04          phys_addr_t end = memblock_end_of_DRAM();
05          int ret;
06
07          if (numa_off)
08                  pr_info("NUMA disabled\n"); /* Forced off on command lin
    e. */
09          pr_info("Faking a node at [mem %#018Lx-%#018Lx]\n", start, end -
    1);
10
11          ret = numa_add_memblk(0, start, end);
12          if (ret) {
13                  pr_err("NUMA init failed\n");
14                  return ret;
15          }
16
17          numa_off = true;
18          return 0;
19  }
```

If there is no NUMA setting, initialize it with a dummy NUMA configuration consisting of one node.

- If the kernel parameter "numa=off" is specified in code lines 7~8, it will output the message "NUMA disabled".
- In line 9 of the code, "Faking a node at..." Outputs a single memory information via message information.
- In lines 11~15 of the code, specify node 0 for both memory memblocks.
- On line 17 of the code, NUMA indicates that it is disabled.

## numa_add_memblk()

arch/arm64/mm/numa.c

```
1  /**
2   * numa_add_memblk() - Set node id to memblk
3   * @nid: NUMA node ID of the new memblk
```

```
 4   * @start: Start address of the new memblk
 5   * @end:  End address of the new memblk
 6   *
 7   * RETURNS:
 8   * 0 on success, -errno on failure.
 9   */
```

```
01  int __init numa_add_memblk(int nid, u64 start, u64 end)
02  {
03          int ret;
04
05          ret = memblock_set_node(start, (end - start), &memblock.memory,
    nid);
06          if (ret < 0) {
07                  pr_err("memblock [0x%llx - 0x%llx] failed to add on node
    %d\n",
08                          start, (end - 1), nid);
09                  return ret;
10          }
11
12          node_set(nid, numa_nodes_parsed);
13          return ret;
14  }
```

Specify the node ID in the memblock area and set the current node in the numa_nodes_parsed bitmap.

## "numa" Early Kernel Parameters

numa_parse_early_param()

arch/arm64/mm/numa.c

```
01  static __init int numa_parse_early_param(char *opt)
02  {
03          if (!opt)
04                  return -EINVAL;
05          if (str_has_prefix(opt, "off"))
06                  numa_off = true;
07
08          return 0;
09  }
10  early_param("numa", numa_parse_early_param);
```

If the "numa=off" kernel parameter is specified, assign true to the global variable numa_off.

# Initializing the NUMA system via the device tree

Refer to the system configuration with four NUMA nodes as follows.

- Each node consists of 4 clusters, for a total of 16 clusters.
- Each cluster has four cores, for a total of 4 cores.
- Since the numa id of the memory node is recognized and passed on by the UEFI firmware, only the information of node 0 is recorded below.

```
#include "hip07.dtsi"

/ {
        model = "Hisilicon Hip07 D05 Development Board";
        compatible = "hisilicon,hip07-d05";

        /* the mem node will be updated by UEFI. */
        memory@0 {
                device_type = "memory";
                reg = <0x0 0x00000000 0x0 0x40000000>;
                numa-node-id = <0>;
        };

        distance-map {
                compatible = "numa-distance-map-v1";
                distance-matrix = <0 0 10>,
                                  <0 1 15>,
                                  <0 2 20>,
                                  <0 3 25>,
                                  <1 0 15>,
                                  <1 1 10>,
                                  <1 2 25>,
                                  <1 3 30>,
                                  <2 0 20>,
                                  <2 1 25>,
                                  <2 2 10>,
                                  <2 3 15>,
                                  <3 0 25>,
                                  <3 1 30>,
                                  <3 2 15>,
                                  <3 3 10>;
```

arch/arm64/boot/dts/hisilicon/hip07.dtsi

```
cpu-map {
        cluster0 {
                core0 {
                        cpu = <&cpu0>;
                };
                core1 {
                        cpu = <&cpu1>;
                };
                core2 {
                        cpu = <&cpu2>;
                };
                core3 {
                        cpu = <&cpu3>;
                };
        };

        ...

        cluster15 {
                core0 {
                        cpu = <&cpu60>;
                };
                core1 {
                        cpu = <&cpu61>;
                };
                core2 {
                        cpu = <&cpu62>;
                };
                core3 {
                        cpu = <&cpu63>;
                };
        };
};

cpu0: cpu@10000 {
        device_type = "cpu";
        compatible = "arm,cortex-a72", "arm,armv8";
        reg = <0x10000>;
        enable-method = "psci";
        next-level-cache = <&cluster0_l2>;
        numa-node-id = <0>;
};

...

cpu63: cpu@70303 {
        device_type = "cpu";
        compatible = "arm,cortex-a72", "arm,armv8";
        reg = <0x70303>;
        enable-method = "psci";
        next-level-cache = <&cluster15_l2>;
        numa-node-id = <3>;
};

cluster0_l2: l2-cache0 {
        compatible = "cache";
```

```
                };

                ...

                cluster15_l2: l2-cache15 {
                        compatible = "cache";
                };
        };
```
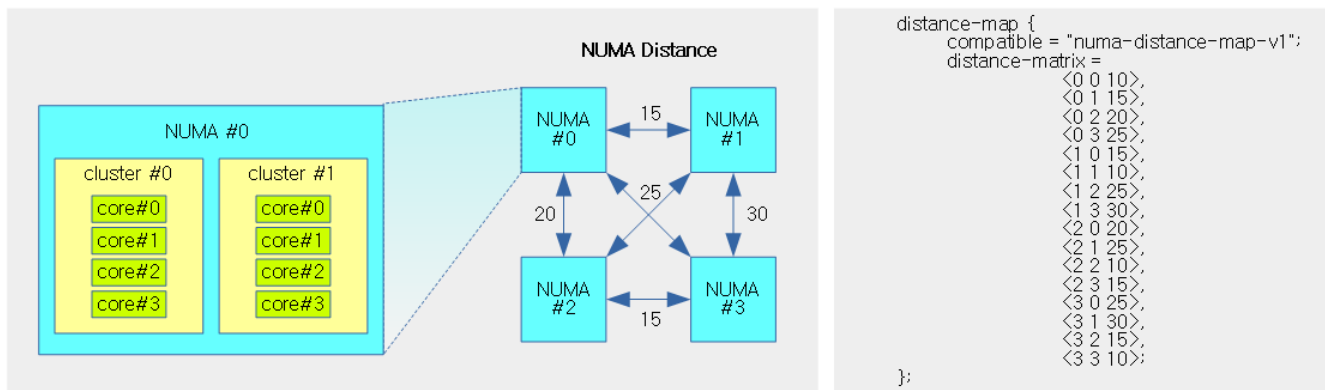
The following figure shows the distances of the NUMA system.



(http://jake.dothome.co.kr/wp-content/uploads/2019/04/numa-distance-1.png)

## of_numa_init()

drivers/of/of_numa.c

```
01  int __init of_numa_init(void)
02  {
03          int r;
04
05          of_numa_parse_cpu_nodes();
06          r = of_numa_parse_memory_nodes();
07          if (r)
08                  return r;
09          return of_numa_parse_distance_map();
10  }
```

It reads the Numa ID from the CPU node and the memory node through the device tree, and parses the distance map.

## of_numa_parse_cpu_nodes()

drivers/of/of_numa.c

```
1  /*
2   * Even though we connect cpus to numa domains later in SMP
3   * init, we need to know the node ids now for all cpus.
4   */
```

```
01  static void __init of_numa_parse_cpu_nodes(void)
02  {
03          u32 nid;
04          int r;
05          struct device_node *np;
```

```
06
07            for_each_of_cpu_node(np) {
08                    r = of_property_read_u32(np, "numa-node-id", &nid);
09                    if (r)
10                            continue;
11
12                    pr_debug("CPU on %u\n", nid);
13                    if (nid >= MAX_NUMNODES)
14                            pr_warn("Node id %u exceeds maximum value\n", ni
   d);
15                    else
16                            node_set(nid, numa_nodes_parsed);
17            }
18 }
```

It traverses the number of CPU nodes and parses them.

- If the "numa-node-id" attribute is greater than or equal to 1, numa_nodes_parsed sets the node in the bitmap.

## of_numa_parse_memory_nodes()

drivers/of/of_numa.c

```
01 static int __init of_numa_parse_memory_nodes(void)
02 {
03        struct device_node *np = NULL;
04        struct resource rsrc;
05        u32 nid;
06        int i, r;
07
08        for_each_node_by_type(np, "memory") {
09                r = of_property_read_u32(np, "numa-node-id", &nid);
10                if (r == -EINVAL)
11                        /*
12                         * property doesn't exist if -EINVAL, continue
13                         * looking for more memory nodes with
14                         * "numa-node-id" property
15                         */
16                        continue;
17
18                if (nid >= MAX_NUMNODES) {
19                        pr_warn("Node id %u exceeds maximum value\n", ni
   d);
20                        r = -EINVAL;
21                }
22
23                for (i = 0; !r && !of_address_to_resource(np, i, &rsrc);
   i++)
24                        r = numa_add_memblk(nid, rsrc.start, rsrc.end +
   1);
25
26                if (!i || r) {
27                        of_node_put(np);
28                        pr_err("bad property in memory node\n");
29                        return r ? : -EINVAL;
30                }
31        }
32
33        return 0;
34 }
```

It traverses the number of memory nodes, parses the nodes, and writes the nodes to the memblock.

- If it is a memory node with the "numa-node-id" attribute, write the node information to the memblock in that area.

## of_numa_parse_distance_map()

drivers/of/of_numa.c

```
01  static int __init of_numa_parse_distance_map(void)
02  {
03          int ret = 0;
04          struct device_node *np;
05
06          np = of_find_compatible_node(NULL, NULL,
07                                       "numa-distance-map-v1");
08          if (np)
09                  ret = of_numa_parse_distance_map_v1(np);
10
11          of_node_put(np);
12          return ret;
13  }
```

It comes from parsing the distance map on nodes with the value of the attribute compatible = "numa-distance-map-v1".

## of_numa_parse_distance_map_v1()

drivers/of/of_numa.c

```
01  static int __init of_numa_parse_distance_map_v1(struct device_node *map)
02  {
03          const __be32 *matrix;
04          int entry_count;
05          int i;
06
07          pr_info("parsing numa-distance-map-v1\n");
08
09          matrix = of_get_property(map, "distance-matrix", NULL);
10          if (!matrix) {
11                  pr_err("No distance-matrix property in distance-map\n");
12                  return -EINVAL;
13          }
14
15          entry_count = of_property_count_u32_elems(map, "distance-matri
    x");
16          if (entry_count <= 0) {
17                  pr_err("Invalid distance-matrix\n");
18                  return -EINVAL;
19          }
20
21          for (i = 0; i + 2 < entry_count; i += 3) {
22                  u32 nodea, nodeb, distance;
23
24                  nodea = of_read_number(matrix, 1);
25                  matrix++;
26                  nodeb = of_read_number(matrix, 1);
27                  matrix++;
28                  distance = of_read_number(matrix, 1);
29                  matrix++;
30
31                  if ((nodea == nodeb && distance != LOCAL_DISTANCE) ||
32                      (nodea != nodeb && distance <= LOCAL_DISTANCE)) {
```

```
33                    pr_err("Invalid distance[node%d -> node%d] = %d
   \n",
34                            nodea, nodeb, distance);
35                    return -EINVAL;
36                }
37
38                numa_set_distance(nodea, nodeb, distance);
39
40                /* Set default distance of node B->A same as A->B */
41                if (nodeb > nodea)
42                        numa_set_distance(nodeb, nodea, distance);
43            }
44
45        return 0;
46  }
```

parse the distance map.

- Read the value of the "distance-matrix" attribute and assign a value to the numa_distance[] array.

## numa_set_distance()

arch/arm64/mm/numa.c

```
01  /**
02   * numa_set_distance() - Set inter node NUMA distance from node to node.
03   * @from: the 'from' node to set distance
04   * @to: the 'to'  node to set distance
05   * @distance: NUMA distance
06   *
07   * Set the distance from node @from to @to to @distance.
08   * If distance table doesn't exist, a warning is printed.
09   *
10   * If @from or @to is higher than the highest known node or lower than z
    ero
11   * or @distance doesn't make sense, the call is ignored.
12   *
13   */
```

```
01  void __init numa_set_distance(int from, int to, int distance)
02  {
03        if (!numa_distance) {
04                pr_warn_once("Warning: distance table not allocated yet
   \n");
05                return;
06        }
07
08        if (from >= numa_distance_cnt || to >= numa_distance_cnt ||
09                        from < 0 || to < 0) {
10                pr_warn_once("Warning: node ids are out of bound, from=%
   d to=%d distance=%d\n",
11                                from, to, distance);
12                return;
13        }
14
15        if ((u8)distance != distance ||
16            (from == to && distance != LOCAL_DISTANCE)) {
17                pr_warn_once("Warning: invalid distance parameter, from
   =%d to=%d distance=%d\n",
18                                from, to, distance);
19                return;
20        }
21
22        numa_distance[from * numa_distance_cnt + to] = distance;
```

```
 23   }
```

numa distace.

- numa_distance specifies a distance value in the [from * numa_distance_cnt + to] array.

## consultation

- NUMA -1- (ARM64 initialization) | Sentence C – Current post
- NUMA -2- (Fallback Node) (http://jake.dothome.co.kr/numa-fallback-node) | Qc
- NODE Bitmap (API) (http://jake.dothome.co.kr/node-api) | Qc


- NUMA binding description
  (https://github.com/torvalds/linux/blob/master/Documentation/devicetree/bindings/numa.txt) |
  kernel.org
- NUMA with Linux (https://lunatine.net/2016/07/14/numa-with-linux/) | Lunatine's Box

---

**LEAVE A COMMENT**

Your email will not be published. Required fields are marked with *

Comments

name *

email *

Website

[ WRITE A COMMENT ]

❮ bootmem_init() – ARM64 (http://jake.dothome.co.kr/bootmem_init-64/)

ARM64 System Key Registers ❯ (http://jake.dothome.co.kr/registers64/)

Munc Blog (2015 ~ 2024)