



# Triton Join Code

For PCI-e

Github repo : <https://github.com/hvdrk/Triton>

Guide : <https://github.com/hvdrk/Triton/asdfn>

# DistributedNumaMem

```
147     let (cache_node, spill_node) =
148         if let MemType::DistributedNumaMem { nodes, .. } = partitions_mem_type {
149             if let [cache_node, spill_node] = *nodes {
150                 Ok((cache_node.node, spill_node.node))
151             } else {
152                 Err(ErrorKind::InvalidArgument(
153                     "Partitioned memory type define exactly two NUMA nodes".to_string(),
154                 ))
155             }
156         } else {
157             Err(ErrorKind::InvalidArgument(
158                 "Partitioned memory type must be DistributedNumaMem".to_string(),
159             ))
160         }?;
```

radix-join/src/execution\_methods/gpu\_triton\_join.rs

Only for DistributedNumaMem!!

```
45 #[derive(Clone, Debug, PartialEq)]
46 pub enum MemType {
47     /// System memory allocated with Rust's global allocator
48     SysMem,
49     /// Aligned system memory allocated with Rust's global allocator
50     ///
51     /// Alignment is specified in bytes.
52     AlignedSysMem { align_bytes: usize },
53     /// NUMA memory allocated on the specified NUMA node and with the specified page type
54     NumaMem { node: u16, page_type: PageType },
55     /// NUMA memory allocated on the specified NUMA node and pinned with CUDA
56     NumaPinnedMem { node: u16, page_type: PageType },
57     /// NUMA memory distributed in proportion to a ratio over multiple NUMA nodes
58     DistributedNumaMem {
59         nodes: Box<[NodeRatio]>,
60         page_type: PageType,
61     },
62     /// NUMA memory distributed over multiple NUMA nodes using a length per node
63     DistributedNumaMemWithLen {
64         nodes: Box<[NodeLen]>,
65         page_type: PageType,
66     },
67     /// CUDA pinned memory (using cudaHostAlloc())
68     CudaPinnedMem,
69     /// CUDA unified memory
70     CudaUniMem,
71     /// CUDA device memory
72     CudaDevMem,
73 }
```

Numa-gpu/src/runtime/allocator.rs

# DistributedNumaMem

## MemType::DistributedNumaMem

```
58 DistributedNumaMem {
59     nodes: Box<[NodeRatio]>,
60     page_type: PageType,
61 },
```

Numa-gpu/src/runtime/allocator.rs

## Memory::DerefMem

```
309 /// A CPU-dereferencable memory type
310 ///
311 /// These memory types can be directly accessed on the host.
312 #[derive(Debug)]
313 pub enum DerefMem<T: DeviceCopy> {
314     /// System memory allocated with Rust's global allocator
315     SysMem(Vec<T>),
316     /// System memory allocated with Rust's global allocator but as a Rust Box
317     ///
318     /// Useful for wrapping aligned memory.
319     BoxedSysMem(Box<[T]>),
320     /// NUMA memory allocated on the specified NUMA node
321     NumaMem(NumaMemory<T>),
322     /// NUMA memory distributed over multiple NUMA nodes
323     DistributedNumaMem(DistributedNumaMemory<T>),
324     /// CUDA pinned memory (using cudaHostAlloc())
325     CudaPinnedMem(LockedBuffer<T>),
326     /// CUDA unified memory
327     CudaUniMem(UnifiedBuffer<T>),
328 }
```

Numa-gpu/src/runtime/memory.rs

## fn alloc\_distributed\_numa

```
356 /// Allocates memory on multiple, specified NUMA nodes.
357 fn alloc_distributed_numa<T: DeviceCopy>(
358     len: usize,
359     nodes: Box<[NodeRatio]>,
360     page_type: PageType,
361 ) -> DerefMem<T> {
362     DerefMem::DistributedNumaMem(DistributedNumaMemory::new_with_ratio(len, nodes, page_type))
363 }
```

Numa-gpu/src/runtime/allocator.rs

## Numa::DistributedNumaMemory<T>

```
394 #[derive(Debug)]
395 pub struct DistributedNumaMemory<T> {
396     ptr: *mut T,
397     len: usize,
398     node_ratios: Box<[NodeRatio]>,
399     page_type: PageType,
400     is_memory_locked: bool,
401     is_page_locked: bool,
402 }
```

Numa-gpu/src/runtime/numa.rs

# DistributedNumaMem

## 1. Allocate memory in VA by mmap()

- void \*mmap(void \*addr, size\_t length, int prot, int flags, int fd, off\_t offset);

prot: protect. Read or write

flags: mapping method. MAP\_ANONYMOUS: not mapped with any file

hugetlb\_flag: pagetype

fd: file descriptor

offset: starting point in file mapping

```
let hugetlb_flags = match page_type {  
  PageType::Huge2MB => libc::MAP_HUGETLB | libc::MAP_HUGE_2MB,  
  PageType::Huge16MB => libc::MAP_HUGETLB | libc::MAP_HUGE_16MB,  
  PageType::Huge1GB => libc::MAP_HUGETLB | libc::MAP_HUGE_1GB,  
  PageType::Huge16GB => libc::MAP_HUGETLB | libc::MAP_HUGE_16GB,  
  PageType::Default | PageType::Small | PageType::TransparentHuge => 0,  
};
```



ptr

```
impl<T> DistributedNumaMemory<T>  
fn new_with_pages
```

```
506 let size = len * size_of::<T>();  
507 let ptr = unsafe {  
508     mmap(  
509         ptr::null_mut(),  
510         size,  
511         libc::PROT_READ | libc::PROT_WRITE,  
512         libc::MAP_PRIVATE | libc::MAP_ANONYMOUS | hugetlb_flags,  
513         0,  
514         0,  
515     )  
516 };  
517
```

Numa-gpu/src/runtime/numa.rs

# DistributedNumaMem

## 2. Set page by madvise()

- int madvise(void \*addr, size\_t length, int advice);  
advice: memory policy

```
impl<T> DistributedNumaMemory<T>  
fn new_with_pages
```

```
523     let advice = match page_type {  
524         PageType::Small => Some(libc::MADV_NOHUGEPAGE),  
525         PageType::TransparentHuge => Some(libc::MADV_HUGEPAGE),  
526         PageType::Default  
527         | PageType::Huge2MB  
528         | PageType::Huge16MB  
529         | PageType::Huge1GB  
530         | PageType::Huge16GB => None,  
531     };  
532  
533     if let Some(advice_flag) = advice {  
534         unsafe {  
535             if madvise(ptr, size, advice_flag) == -1 {  
536                 let err = IoError::last_os_error();  
537                 std::result::Result::Err::<(), _>(err).expect("Failed to madvise memory");  
538             }  
539         }  
540     }
```

Numa-gpu/src/runtime/numa.rs



ptr

# DistributedNumaMem

## 3. Binding memory with each Numa node by mbind()

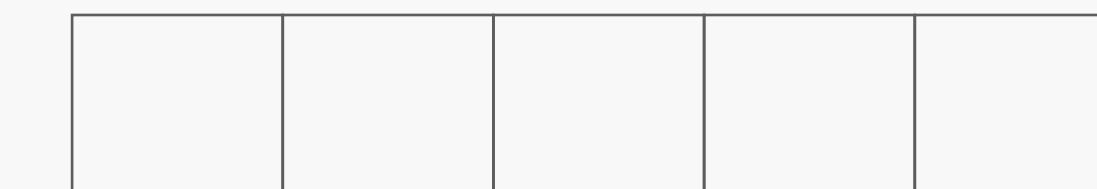
- int mbind(void \*start, unsigned long len, int mode, const unsigned long \*nodemask, unsigned long maxnode, unsigned int flags);

mode: memory policy. BIND->allocate

nodemask: node bitmask. Ex) 101 -> bind node 0,2

- In code, ptr ~ ptr+slice -> cache node -> device numa mem  
ptr +slice ~ -> spill node -> host numa mem

Bind with node255 at iter 0 (device)	Bind with node0 at iter 1 (host)



ptr

impl<T> DistributedNumaMemory<T>  
fn new\_with\_pages

```
553 let final_node_ratios = node_pages
554 .iter()
555 .scan(
556     0,
557     |page_offset,
558     &NodeLen {
559         node,
560         len: page_len,
561     }| {
562         let old = *page_offset;
563         *page_offset = *page_offset + page_len;
564         Some((node, old, page_len))
565     },
566 )
567 .map(|(node, page_offset, page_len)| {
568     let mut node_set = CpuSet::new();
569     node_set.add(node);
570
571     unsafe {
572         let slice = slice::from_raw_parts(
573             ptr.add(page_offset as usize * page_size),
574             page_len as usize * page_size,
575         );
576
577         // Note that mbind fails with `EINVAL` for HugeTLB mappings if `size`
578         // isn't a multiple of the page size
579         mbind(slice, MemPolicyModes::BIND, node_set, MemBindFlags::STRICT?);
580     }
581
582     Ok(NodeRatio {
583         node,
584         ratio: Ratio::<usize>::new(page_len, pages),
585     })
586 })
587 .collect::<Result<Box<[NodeRatio]>>>()
588 .expect("Failed to mbind memory to the specified NUMA nodes");
```

Numa-gpu/src/runtime/numa.rs



# PCI-e

## Problem

- In PCI-e, device memory doesn't be caught as Numa node.
  - > but DistributedNumaMemory need two Numa node.
    - One for device, the other for host



# PCI-e

```
425 let mut inner_rel_partitions = PartitionedRelation::new(
426     data.build_relation_key.len(),
427     histogram_algorithmfst.either(|cpu| cpu.into(), |gpu| gpu.into()),
428     radix_bits.pass_radix_bits(RadixPass::First).unwrap(),
429     max_chunks_1st,
430     inner_rel_alloc(cache_bytes_inner / mem::size_of::<Tuple<T, T>>()), //distributednuma allocator
431     Allocator::mem_alloc_fn(offsets_mem_type.clone()), //numa allocator
432 );
433
```

radix-join/src/execution\_methods/gpu\_triton\_join.rs

```
460 radix_pnr.partition(
461     RadixPass::First,
462     data.build_relation_key.as_launchable_slice(),
463     data.build_relation_payload.as_launchable_slice(),
464     &mut inner_rel_partition_offsets,
465     &mut inner_rel_partitions,
466     &stream,
467 );
```

```
188 let mut radix_pnr = GpuRadixPartitioner::new(
189     histogram_algorithmfst.gpu_or_else(|cpu_algo| cpu_algo.into()),
190     partition_algorithmfst,
191     radix_bits.clone(),
192     grid_size,
193     block_size,
194     dmembuffer_bytes,
195 );
196 radix_pnr.preallocate_partition_state::<T>(RadixPass::First)?;
```

```
323 #[derive(Debug)]
324 pub struct PartitionedRelation<T: DeviceCopy> {
325     pub relation: Mem<T>,
326     pub offsets: Mem<u64>,
327     len: usize,
328     chunks: u32,
329     radix_bits: u32,
330 }
331
332 impl<T: DeviceCopy> PartitionedRelation<T> {
333     /// Creates a new partitioned relation, and automatically includes the
334     /// necessary padding and metadata.
335     pub fn new(
336         len: usize,
337         histogram_algorithm_type: HistogramAlgorithmType,
338         radix_bits: u32,
339         max_chunks: u32,
340         partition_alloc_fn: MemAllocFn<T>,
341         offsets_alloc_fn: MemAllocFn<u64>,
342     ) -> Self {
343         let chunks: u32 = match histogram_algorithm_type {
344             HistogramAlgorithmType::Chunked => max_chunks,
345             HistogramAlgorithmType::Contiguous => 1,
346         };
347
348         let padding_len = padding_len::<T>();
349         let num_partitions = fanout(radix_bits) as usize;
350         let relation_len = len + (num_partitions * chunks as usize) * padding_len as usize;
351
352         let relation = partition_alloc_fn(relation_len);
353         let offsets = offsets_alloc_fn(num_partitions * chunks as usize);
354
355         Self {
356             relation,
357             offsets,
358             chunks,
359             radix_bits,
360             len,
361         }
362     }
363 }
```

sql-ops/src/partition/partitioned\_relation.rs

```
314 #[derive(Debug)]
315 pub struct GpuRadixPartitioner {
316     radix_bits: RadixBits,
317     prefix_sum_algorithm: GpuHistogramAlgorithm,
318     partition_algorithm: GpuRadixPartitionAlgorithm,
319     prefix_sum_state: PrefixSumState,
320     partition_state: RadixPartitionState,
321     grid_size: GridSize,
322     block_size: BlockSize,
323     rp_block_size: BlockSize,
324     dmembuffer_bytes: usize,
325 }
```

```
536 pub fn partition<T: DeviceCopy + GpuRadixPartitionable>(
537     &mut self,
538     pass: RadixPass,
539     partition_attr: LaunchableSlice<'_, T>,
540     payload_attr: LaunchableSlice<'_, T>,
541     partition_offsets: &mut PartitionOffsets<Tuple<T, T>>,
542     partitioned_relation: &mut PartitionedRelation<Tuple<T, T>>,
543     stream: &Stream,
544 ) -> Result<()> {
545     T::partition_impl(
546         self,
547         pass,
548         partition_attr,
549         payload_attr,
550         partition_offsets,
551         partitioned_relation,
552         stream,
553     )
554 }
```

```
1178 GpuRadixPartitionAlgorithm::SSWMCv2 => {
1179     let name = std::ffi::CString::new(
1180         stringify!([<gpu_chunked_sswmc_radix_partition_v2_ $Suffix _ $Suffix>])
1181     ).unwrap();
1182     let mut function = module.get_function(&name)?;
1183     function.set_max_dynamic_shared_size_bytes(max_shared_mem_bytes)?;
1184
1185     unsafe {
1186         launch!(
1187             function<<<
1188                 grid_size,
1189                 rp_block_size,
1190                 max_shared_mem_bytes,
1191                 stream
1192             >>>(
1193                 args.clone(),
1194                 max_shared_mem_bytes
1195             ));
1196     }
1197 }
```

Function in

sql-ops/cudautils/radix\_partition.cu



# PCI-e

```
688 template <typename K, typename V>
689 __device__ void gpu_chunked_sswc_radix_partition_v2(
690     RadixPartitionArgs &args, uint32_t shared_mem_bytes) {
```

Sql-ops/cudautils/radix\_partition.cu

```
101 struct RadixPartitionArgs {
102     // Inputs
103     const void *const __restrict__ join_attr_data;
104     const void *const __restrict__ payload_attr_data;
105     std::size_t const data_length;
106     uint32_t const padding_length;
107     uint32_t const radix_bits;
108     uint32_t const ignore_bits;
109     const unsigned long long *const __restrict__ partition_offsets;
110
111     // State
112     uint32_t *const __restrict__ tmp_partition_offsets;
113     char *const __restrict__ l2_cache_buffers;
114     char *const __restrict__ device_memory_buffers;
115     uint64_t const device_memory_buffer_bytes;
116
117     // Outputs
118     void *const __restrict__ partitioned_relation;
119 };
```

Sql-ops/include/gpu\_radix\_partition.h

- Just use partitioned\_relation in continuous VA
- Except for the NUMA node approach, different buffers can't be mapped into a contiguous VA



# Solution

- 1) use UnifiedBuffer
- 2) Adaptive memory management
- 3) use custom Buffer

} Change only rust code

Change almost every code

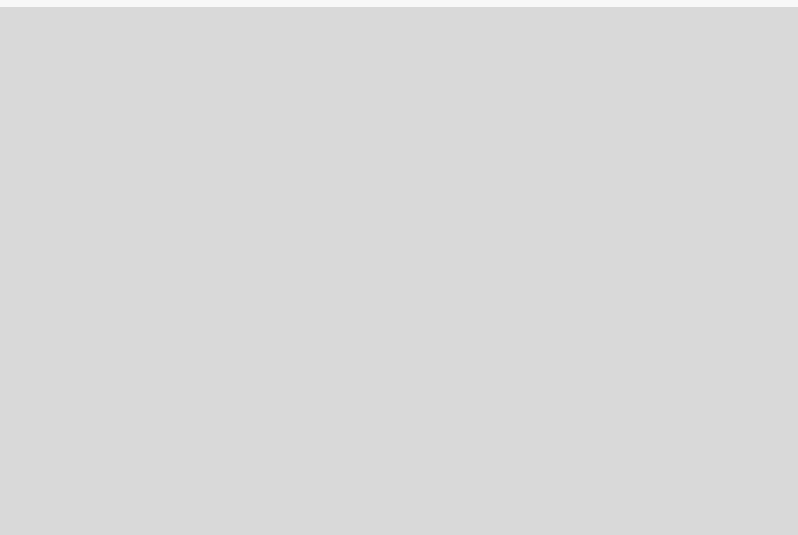
# Solution 1. UnifiedBuffer

rustcuda::memory::UnifiedBuffer

- <https://bheisler.github.io/RustaCUDA/rustacuda/memory/struct.UnifiedBuffer.html>
- Use Nvidia's API cuMemAllocManaged()
- Works well ..., but **extremely slow** when exceeding device's memory
- Can't map manually, and it is assumed that a page fault occurs with every one access for exceeded



## Solution 2. Adaptive memory management

- Use only host memory when exceed device's memory capacity.
  - When the data size exceeds the device memory capacity, there is a slowdown
  - However, when the data size is sufficiently large, the performance difference becomes negligible.
- 



## Solution 3. Custom Buffer

- The method is quite simple, but it remains **unimplemented** due to the large number of required modifications
- Make new data struct include device buffer and host numa memory
- Use `.get()` , `.set()`, ... and adjust for every file include cuda file(.cu), header file(.h), cpp file(.cpp)



# Solution 2. Changes



# radix-join/Cargo.toml

Add cust = "0.3.2" in [dependencies]

- cust : use CUDA in rust
- For get device's free memory

```
radix-join > Cargo.toml
 9  edition = "2018"
10
11  [dependencies]
12  cstr = "0.2.8"
13  csv = "~1.1.1"
14  hostname = "~0.1.5"
15  itertools = "0.9"
16  num-rational = "~0.2.0"
17  num-traits = "~0.2.0"
18  rayon = "~1.2.0"
19  rustacuda = { git = "https://github.com/LutzCle/RustaCUDA", branch = "custom_mods_10_2" }
20  serde = "~1.0.76"
21  serde_derive = "~1.0.76"
22  serde_repr = "~0.1"
23  structopt = "0.3"
24  cust = "0.3.2"
```

# radix-join/src/main.rs

## Change `fn set_partitions_mem`

- Triton join arise error when memory type of 1st partitioned\_relation is not DistributedNuma
- Change to NumaPinned when PCI-e

## `fn set_partitions_mem`

```
348 impl CmdOpt {
358     fn set_partitions_mem(
359         &mut self,
360         cache_location: Option<u16>,
361         overflow_location: u16,
362     ) -> Result<> {
363         if self.execution_method == ArgExecutionMethod::GpuTritonJoinTwoPass {
364             if cache_location != None {
365                 let cache_location = cache_location.ok_or_else(|| {
366                     ErrorKind::RuntimeError(
367                         "Failed to set the cache NUMA location. Are you using PCI-e?".to_string(),
368                     )
369                 })?;
370
371                 if ArgMemType::DistributedNuma != self.partitions_mem_type {
372                     self.partitions_mem_type = ArgMemType::DistributedNuma;
373                     self.partitions_location = vec![cache_location, overflow_location];
374                     self.partitions_proportions = vec![0, 0];
375                 } else if self.partitions_location.len() != 2 {
376                     let e = format!(
377                         "Invalid argument: --partitions-location must specify \
378                         exactly two locations when combined with --execution-method \
379                         GpuTritonJoin\n\
380                         The default locations are: --partitions-location {},{}\n",
381                         cache_location, overflow_location
382                     );
383                     Err(ErrorKind::InvalidArgument(e));
384                 }
385             }
386             else {
387                 self.partitions_mem_type = ArgMemType::NumaPinned;
388                 self.partitions_location = vec![0, 0];
389                 self.partitions_proportions = vec![0, 0];
390             }
391         }
392
393         ok(())
394     }
395 }
```

# radix-join/src/gpu\_triton\_join.rs

## Add library

- `sql_ops::partition::{fanout, HistogramAlgorithmType},`  
`sql_ops::partition::partitioned_relation::padding_len`  
-> compute 1st partitioned\_relation's total size
- `cust::memory::mem_get_info`  
-> get device's free memory

# radix-join/src/gpu\_triton\_join.rs

## Change `fn gpu_triton_join`

- Change parameter `partitions_mem_type` to `mut`  
: can be changed after
- Set `cache_node`, `spill_node` to 0(not used) when not `DistributedNumaMem`

```
102 pub fn gpu_triton_join<T>(  
103     data: &mut JoinData<T>,   
104     hashing_scheme: HashingScheme,   
105     histogram_algorithm_fst: DeviceType<CpuHistogramAlgorithm, GpuHistogramAlgorithm>,   
106     histogram_algorithm_snd: DeviceType<CpuHistogramAlgorithm, GpuHistogramAlgorithm>,   
107     partition_algorithm_fst: DeviceType<CpuRadixPartitionAlgorithm, GpuRadixPartitionAlgorithm>,   
108     partition_algorithm_snd: DeviceType<CpuRadixPartitionAlgorithm, GpuRadixPartitionAlgorithm>,   
109     radix_bits: &RadixBits,   
110     dmem_buffer_bytes: usize,   
111     max_partitions_cache_bytes: Option<usize>,   
112     threads: usize,   
113     cpu_affinity: CpuAffinity,   
114     mut partitions_mem_type: MemType,   
115     stream_state_mem_type: MemType,   
116     page_type: PageType,   
117     partition_dim: (&GridSize, &BlockSize),   
118     join_dim: (&GridSize, &BlockSize),   
119 ) -> Result<i64, RadixJoinPoint>
```

```
150 let (cache_node, spill_node) =   
151     if let MemType::DistributedNumaMem { nodes, .. } = partitions_mem_type.clone() {   
152         if let [cache_node, spill_node] = *nodes {   
153             Ok((cache_node.node, spill_node.node))   
154         } else {   
155             Err(ErrorKind::InvalidArgument(   
156                 "Partitioned memory type define exactly two NUMA nodes".to_string(),   
157             ))   
158         }   
159     } else {   
160         Ok((0,0))   
161     }?;
```

# radix-join/src/gpu\_triton\_join.rs

## Change `fn gpu_triton_join`

- Change `offsets_mem_type` to `NumaPinnedMem` when not `DistributedNumaMem`  
: device can't access to `NumaMem`

```
183 |         let offsets_mem_type = match partitions_mem_type {  
184 |             MemType::DistributedNumaMem => MemType::NumaMem {  
185 |                 node: spill_node,  
186 |                 page_type,  
187 |             },  
188 |             _ => partitions_mem_type.clone()  
189 |         }
```

- Get device's free memory

```
420 |         // get device's free memory  
421 |         let free_mem = if let Ok((free_mem, _)) = mem_get_info() {  
422 |             free_mem - GPU_MEM_SLACK_BYTES  
423 |         } else {  
424 |             0  
425 |         };
```

# radix-join/src/gpu\_triton\_join.rs

## Change `fn gpu_triton_join`

- Compute 1st partitioned relation's total size

```
433 // compute 1st partitioned relation's total size
434 let padding_len = padding_len::<Tuple<T, T>>();
435 let chunks: u32 = match histogram_algorithmfst.either(|cpu| cpu.into(), |gpu| gpu.into()) {
436     HistogramAlgorithmType::Chunked => max_chunks_1st,
437     HistogramAlgorithmType::Contiguous => 1,
438 };
439 let num_partitions = fanout(radix_bits.pass_radix_bits(RadixPass::First).unwrap()) as usize;
440 let inner_relation_len = data.build_relation_key.len() + (num_partitions * chunks as usize) * padding_len as usize;
441 let inner_relation_size = inner_relation_len * (mem::size_of::<Tuple<T, T>>() as usize);
442 let outer_relation_len = data.probe_relation_key.len() + (num_partitions * chunks as usize) * padding_len as usize;
443 let outer_relation_size = outer_relation_len * (mem::size_of::<Tuple<T, T>>() as usize);
444 let total_relation_size = inner_relation_size + outer_relation_size;
```

- Change `partitions_mem_type` when device's free memory is larger than total relation size

```
449 // if device's free memory is larger than total_relation_size, use device's memory
450 if free_mem > total_relation_size {
451     partitions_mem_type = MemType::CudaDevMem;
452 }
```



# radix-join/src/gpu\_triton\_join.rs

## Change `fn gpu_triton_join`

- Revise inner/outer relation allocator to not spill when not DistributedNumaMem

```
466 let (inner_rel_alloc, cached_build_tuples) = match partitions_mem_type {
467     MemType::DistributedNumaMem {..} => Allocator::mem_spill_alloc_fn(CacheSpillType::CacheAndSpill {
468         cache_node,
469         spill_node,
470         page_type,
471     }),
472     _ => Allocator::mem_spill_alloc_fn(CacheSpillType::NoSpill(partitions_mem_type.clone()))
473 };
474
```

```
484 let (outer_rel_alloc, cached_probe_tuples) = match partitions_mem_type {
485     MemType::DistributedNumaMem {..} => Allocator::mem_spill_alloc_fn(CacheSpillType::CacheAndSpill {
486         cache_node,
487         spill_node,
488         page_type,
489     }),
490     _ => Allocator::mem_spill_alloc_fn(CacheSpillType::NoSpill(partitions_mem_type.clone()))
491 };

```

# sql-ops/src/

## Change private fn to public

- Used for compute 1st partitioned relation's total size

```
21 pub mod cpu_radix_partition;  
22 pub mod gpu_radix_partition;  
23 mod partition_input_chunk;  
24 pub mod partitioned_relation;  
  
50 pub fn fanout(radix_bits: u32) -> u32 {  
51     1 << radix_bits  
52 }
```

sql-ops/src/partitions.rs

```
26 /// Convert padding bytes into padding length for the type `T`  
27 pub fn padding_len<T: Sized>() -> u32 {  
28     crate::constants::PADDING_BYTES / mem::size_of::<T>() as u32  
29 }
```

sql-ops/src/partition/partitioned\_relation.rs