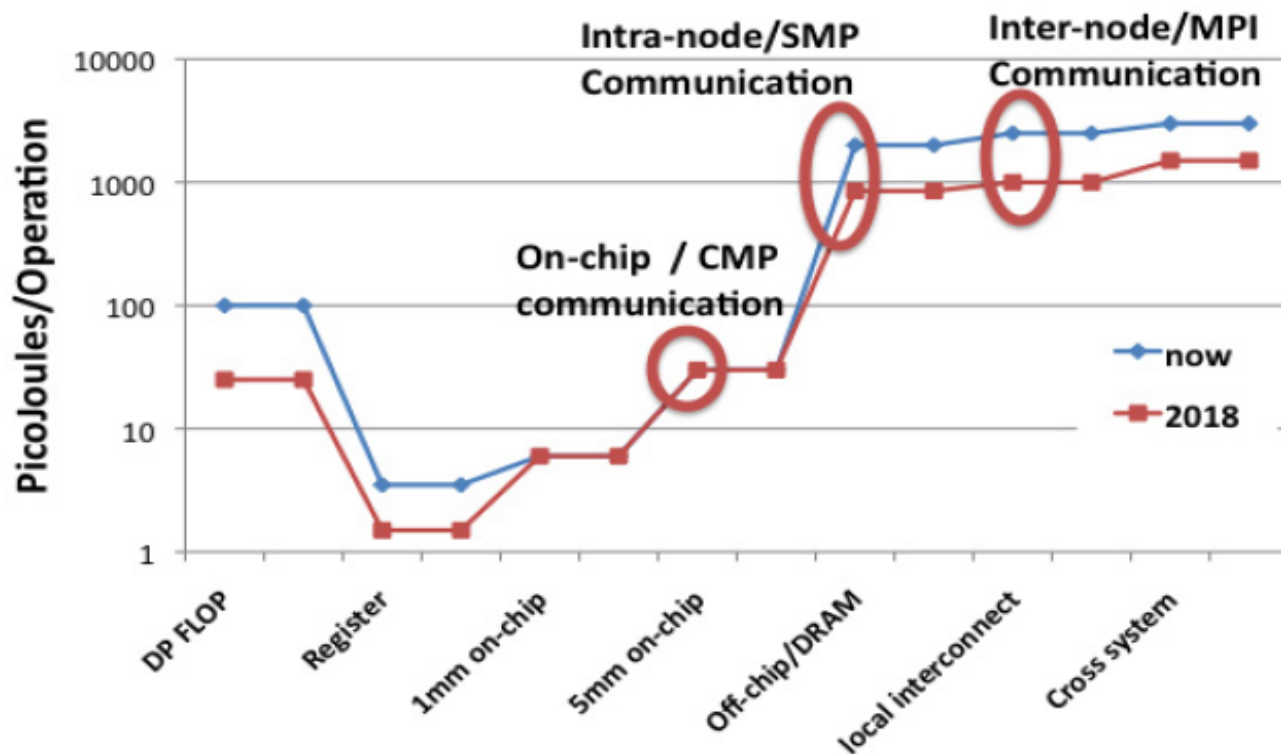# CPU自旋锁优化

2019年10月14日

蚂蚁金服
ANT FINANCIAL

# 内容

- 数据迁移问题 & 对性能影响
- Spinlock Issues on Multi-Socket systems
    MCS Spinlock

    Critical Section Integration

    NUMA Aware Spinlock
- NUMA Spinlock
- NUMA Spinlock Example
- NUMA Spinlock Result
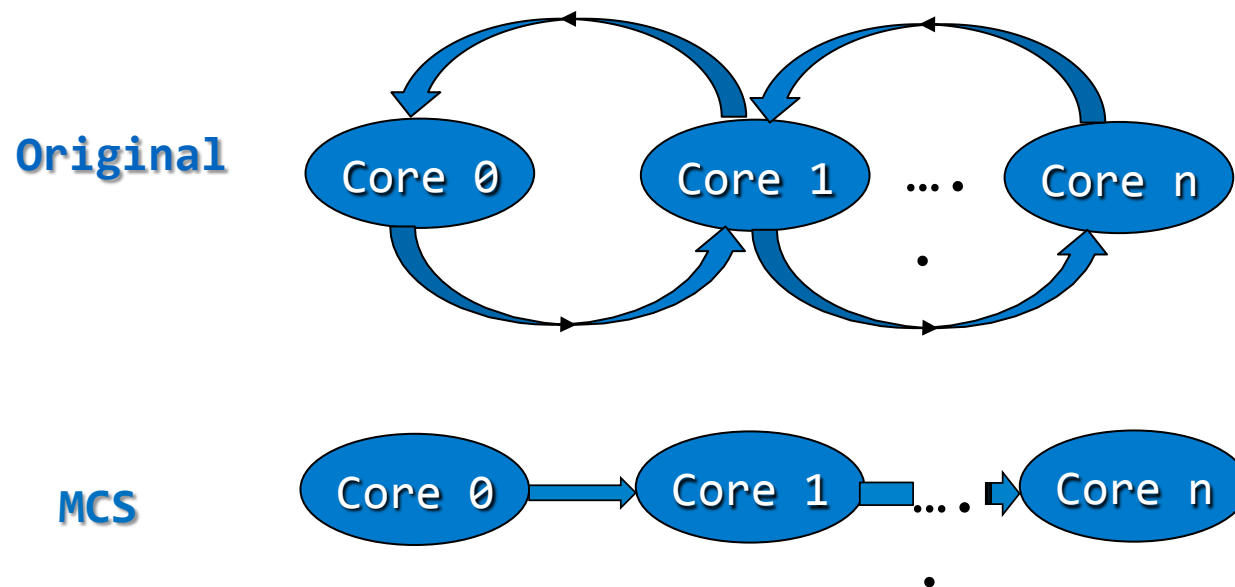
数据迁移是体系架构的主要瓶颈



双浮运算需要 10PJ
数据移动需要 1000PJ

- On multi-socket multiple core systems,
  Access Latency:

  remote socket > local socket > sibling core >
  local core

# MCS Spinlock



Original

Core 0     Core 1    ....    Core n

MCS

Core 0 → Core 1 → .... → Core n

**MCS helps us to reduce useless lock movement in spinning stage.**

# Critical Section Integration(CSI)

## Core 0

```
void * work(void *V)
{
   v->result = CS(v->A);}
   Acquire LOCK
   If(contention) {
      send work
      while(result);
   }else {
      CS(A);
      If(incoming request)
         run work list
      Release LOCK
   }
   print result
```

## Core n

```
void * work(void *V)
{
   v->result = CS(v->A);}
   Acquire LOCK
   If(contention) {
      send work
      while(result);
   }else {
      CS(A);
      If(incoming request)
         run work list
      Release LOCK
   }
   print result
```

```
LOCK X
result =
CS(A);
UNLOCK X
print result
```

**CSI helps us to reduce share data movement.**

# NUMA Aware Spinlock(NAS)
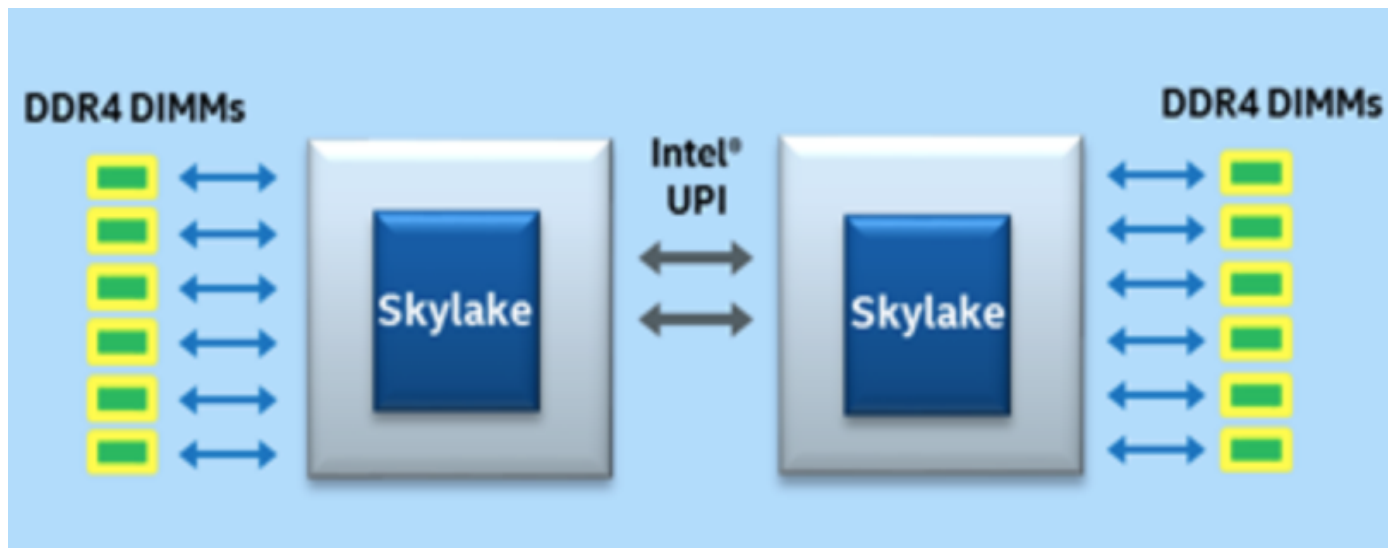
```
Acquire LOCK from Socket 0/1
   if (contention)
      append work list & waiting

Acquire Global LOCK
   if (contention)
      waiting

Run local work list

Release LOCK from Socket 0/1

Release Global LOCK
```
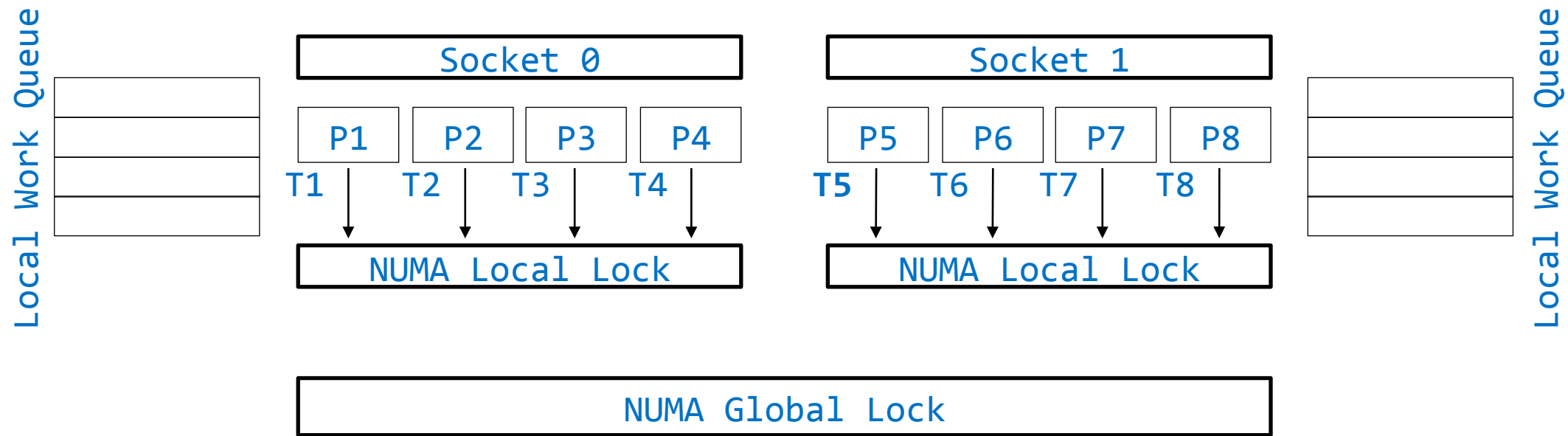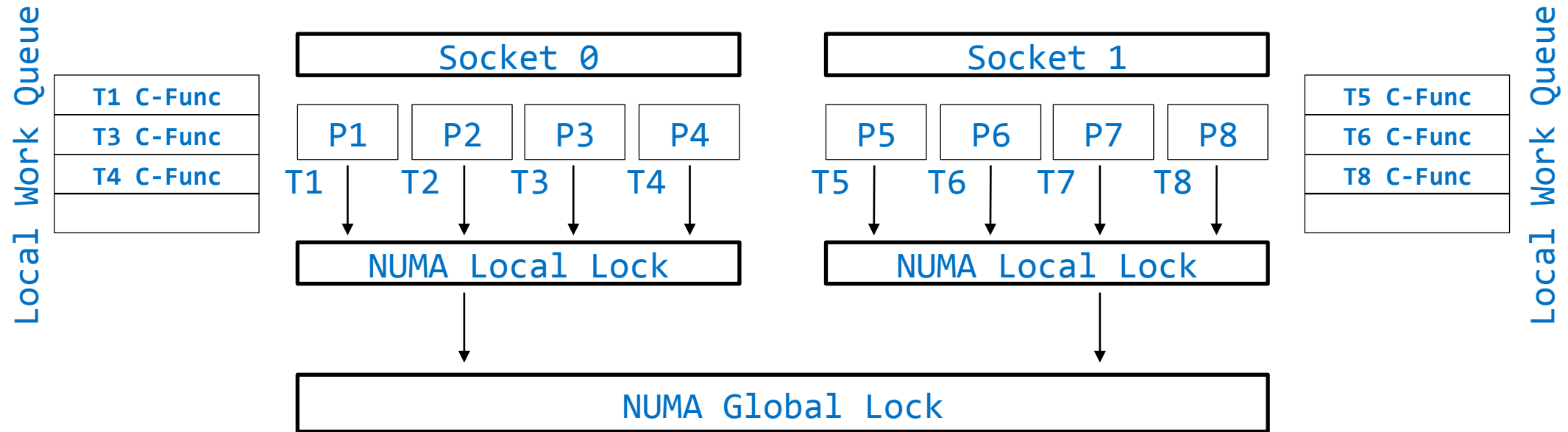
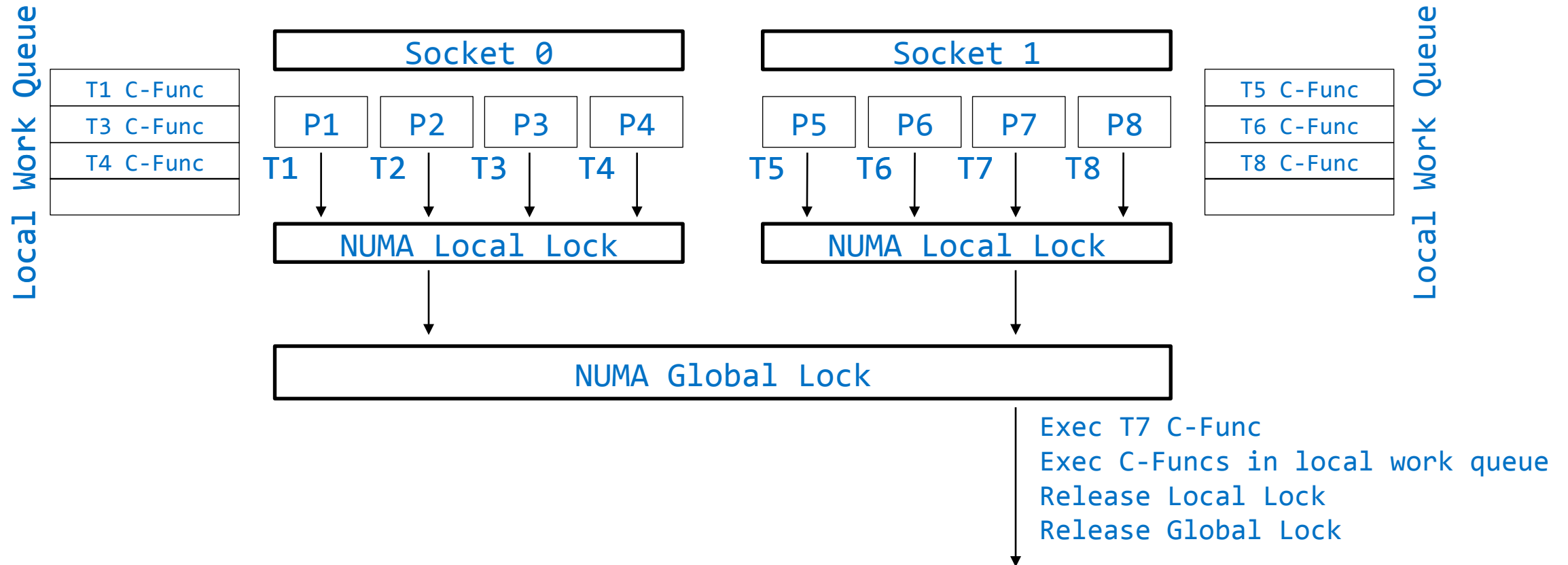**NAS helps us to reduce off-chip lock movement.**

# NUMA Spinlock (2)

# NUMA Spinlock (3)



Local Work Queue

| |
|---|
| T1 C-Func |
| T3 C-Func |
| T4 C-Func |
| |

Socket 0

| P1 | P2 | P3 | P4 |
|---|---|---|---|

T1    T2    T3    T4

NUMA Local Lock

Socket 1

| P5 | P6 | P7 | P8 |
|---|---|---|---|

T5    T6    T7    T8

NUMA Local Lock

Local Work Queue

| |
|---|
| T5 C-Func |
| T6 C-Func |
| T8 C-Func |
| |

NUMA Global Lock

Exec T7 C-Func
Exec C-Funcs in local work queue
Release Local Lock
Release Global Lock

# NUMA Spinlock (5)



Local Work Queue

| T1 C-Func |
|-----------|
| T3 C-Func |
| T4 C-Func |
|           |

**Socket 0**

| P1 | P2 | P3 | P4 |
|----|----|----|----|

T1    T2    T3    T4

**NUMA Local Lock**

**Socket 1**

| P5 | P6 | P7 | P8 |
|----|----|----|----|

T5    T6    T7    T8

Local Work Queue

**NUMA Global Lock**

Exec T2 C-Func
Exec C-Funcs in local work queue
Release Local Lock
Release Global Lock

Local Work Queue

| Socket 0 | | | |
|---|---|---|---|
| P1 | P2 | P3 | P4 |

T1    T2    T3    T4

| Socket 1 | | | |
|---|---|---|---|
| P5 | P6 | P7 | P8 |

T5    T6    T7    T8

Local Work Queue

All works on socket 0,1 are completed.

**NUMA spinlock helps us to minimize cross-socket traffic as well as localize the serialized workload to one core for execution.**

# NUMA Spinlock Example

```
struct numa_spinlock *lock;

work_thread (void *arg) {
    struct work_todo_argument
work_todo_arg;
    struct numa_spinlock_info lock_info;
    if (numa_spinlock_init (lock,
&lock_info)) {
        printf ("numa_spinlock_init
failure: %m\n");
        exit (1);
    }
    work_todo_arg.arg = arg;
    lock_info.argument = &work_todo_arg;
    lock_info.workload = work_todo;
    numa_spinlock_apply (&lock_info);
    return lock_info.result;
}
```

A 'numa_spinlock' pointer uniquely identifies a NUMA spinlock object.

This data type uniquely identifies a NUMA spinlock information object for a thread.

Initialize the NUMA spinlock information block pointed to by INFO with a NUMA spinlock pointer LOCK. The return value is '0' on success and '-1' on failure.

A pointer to argument passed to the WORKLOAD function pointer.

A function pointer to the workload function serialized by spinlock.

Apply for spinlock with a NUMA spinlock information block pointed to by INFO.

When 'numa_spinlock_apply' returns, the spinlock is released and the RESULT member of INFO contains the return value of the WORKLOAD member.

# NUMA Spinlock Example (2)

- static void *work_todo (void *v) {
-     struct work_todo_argument *p = v;
-     void *ret_val;
-     ret_val = serialized_work (p->arg);
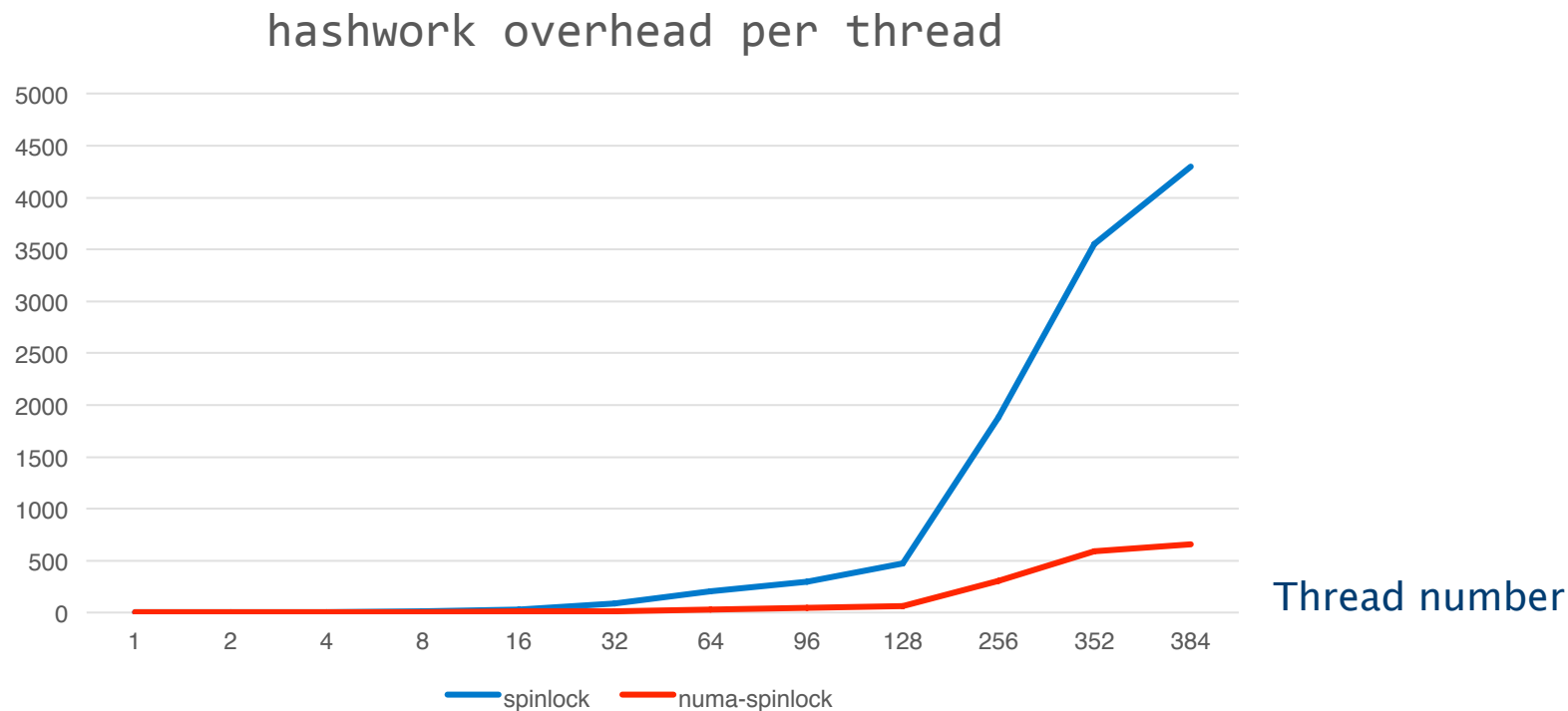-     return ret_val;
- }

> Do the real work with p->arg.

> Return value is set to lock_info.result.

```
static void __attribute__((constructor)) init_numaspinlock (void) {
    lock = numa_spinlock_alloc ();
}

static void __attribute__((destructor)) destroy_numaspinlock (void) {
    numa_spinlock_free (lock);
}
```

# Micro-Benchmark Result

## hashwork overhead per thread



已经提交硬件优化，解决最后的问题，基本达到理论值，自旋锁问题即将解决

1%的串行化导致 64核心只有 40核心的结果　Speed up = T / ((0.99T/64) + 0.01T) = 39.2倍

因此建议尽量不要使用同步机制，不用性能才最好 😁

THANK YOU!