# perf

(Linux profiling with performance counter)

**Multicore Programming** 



### Introduction

• What is Perf?

Perf Commands

Example



### What is Perf?

• 리눅스의 시스템 성능 측정 도구

• Kernel, User mode 모두 측정 가능

• 여러가지 종류의 Event(Hardware, Software, Hardware cache, ...)에 대한 측정 가능

• 최신 리눅스 커널에 포함



### **Perf Commands**

- perf list perf를 통해 측정 가능한 event의 종류를 확인한다.
- perf stat

  Event count를 얻는다.
- perf record

  Events를 sampling하여 output 파일을 생성한다.
- perf report perf record를 통해 생성된 file을 분석한다.
- perf annotate perf record를 통해 생성된 file을 instruction level에서 분석한다.
- more Commands, but not today



### Perf Commands – perf list

#### perf를 통해 측정 가능한 event의 종류를 확인한다.

#### \$ sudo perf list

```
mrbin2002@ubuntu:~/TA multicore/prac mutex$ sudo perf list
List of pre-defined events (to be used in -e):
  cpu-clock
                                                       [Software event]
  task-clock
                                                       [Software event]
  page-faults OR faults
                                                       [Software event]
  context-switches OR cs
                                                       [Software event]
  cpu-migrations OR migrations
                                                       [Software event]
  minor-faults
                                                       [Software event]
  major-faults
                                                       [Software event]
  alignment-faults
                                                       [Software event]
  emulation-faults
                                                       [Software event]
                                                      [Software event]
  dummy
  power/energy-cores/
                                                       [Kernel PMU event]
  power/energy-gpu/
                                                       [Kernel PMU event]
  power/energy-pkg/
                                                       [Kernel PMU event]
                                                       [Raw hardware event descriptor]
  LNNN
  cpu/t1=v1[,t2=v2,t3 ...]/modifier
                                                       [Raw hardware event descriptor]
   (see 'man perf-list' on how to encode it)
                                                      [Hardware breakpoint]
 mem:<addr>[:access]
```



### Perf Commands – perf stat

타겟 프로세스가 실행되는 동안 발생하는 event의 count를 측정한다.

#### \$ sudo perf stat ./a.out

```
mrbin2002@ubuntu:~/TA_multicore/prac_mutex$ sudo perf stat ./a.out
thread 30770944, local count: 1250000
thread 22378240, local count: 1250000
thread 13985536, local count: 1250000
thread 5592832, local count: 1250000
thread -2799872, local count: 1250000
thread -11192576, local count: 1250000
thread -19585280, local count: 1250000
thread -27977984, local count: 1250000
thread -36370688, local count: 1250000
global count: 11250000
Performance counter stats for './a.out':
      11574.015841
                       task-clock (msec)
                                                # 6.143 CPUs utilized
          358,797
                       context-switches
                                                # 0.031 M/sec
                       cpu-migrations
                                                # 0.502 K/sec
            5.814
                                                # 0.007 K/sec
                       page-faults
                       cvcles
                                                     0.000 GHz
                       stalled-cycles-frontend #
                                                     0.00% frontend cycles idle
                       stalled-cycles-backend
                                                     0.00% backend cycles idle
                       instructions
                       branches
                                                     0.000 K/sec
                       branch-misses
                                                     0.000 K/sec
      1.884165349 seconds time elapsed
```



## Perf Commands – perf record

#### Event를 sampling하여 output file을 생성한다.

#### \$ sudo perf record -g ./a.out

```
mrbin2002@ubuntu:~/TA_multicore/prac_mutex$ sudo perf record -g ./a.out thread -554502400, local count: 1250000 thread -562895104, local count: 1250000 thread -571287808, local count: 1250000 thread -579680512, local count: 1250000 thread -588073216, local count: 1250000 thread -596465920, local count: 1250000 thread -604858624, local count: 1250000 thread -613251328, local count: 1250000 thread -621644032, local count: 1250000 global count: 11250000 [ perf record: Woken up 8 times to write data ] [ perf record: Captured and wrote 3.018 MB perf.data (~131878 samples) ]
```

(process를 실행시키고 종료될 때 까지의 event sampling)



# Perf Commands – perf record

#### Event를 sampling하여 output file을 생성한다.

#### \$ sudo perf record -g -p 24123

```
mrbin2002@ubuntu:~/TA_multicore/prac_mutex$ ps u
USER
          PID %CPU %MEM
                         VSZ
                               RSS TTY
                                           STAT START
                                                       TIME COMMAND
mrbin20+ 4962 0.0 0.1 29288
                                                       0:00 /bin/bash
                              6788 pts/7
                                           Ss
                                               Sep21
5:40 vi projec
                                               07:23
                                                       0:00 /bin/bash
mrbin20+ 11292  0.0  0.1  29308
                              6980 pts/7
                                                07:25
mrbin20+ 15421 0.0 0.1 29276 6776 pts/8
                                         Ss+ Sep21
                                                       0:00 /bin/bash
mrbin20+ <u>16771</u> 0.0 0.1 29332 7000 pts/7
                                                Sep21
                                                       0:00 /bin/bash
mrbin20+ 24123 758 0.0 826096 1996 pts/7
                                           Sl+ 23:11 1:53 ./a.out
mrbin20+ 24271 0.0 0.0 23868 2540 pts/0
                                           R+
                                               23:11
                                                       0:00 ps u
mrbin20+ 31826  0.0  0.1  29276  6716 pts/0
                                           Ss
                                               Sep21
                                                       0:00 bash
mrbin20+ 32161  0.6  0.2  58428  11112  pts/7
                                               Sep21 14:19 vi project
mrbin2002@ubuntu:~/TA multicore/prac mutex$ sudo perf record -g -p 24123
[sudo] password for mrbin2002:
 perf record: Woken up 115 times to write data ]
 perf record: Captured and wrote 19.230 MB perf.data (~840175 samples) ]
```

(이미 실행중인 process에 대한 event sampling)



# Perf Commands – perf report

#### perf record를 통해 생성된 file을 분석한다.

\$ sudo perf record -g graph --no-children

```
Samples: 32K of event 'cpu-clock', Event count (approx.): 8086500000
 Overhead Command Shared Object
                                        Symbol
    74.04% a.out
                    [kernel.kallsyms]
                                        [k] raw spin lock
   - raw spin lock
      - 42.96% do_futex
          sys futex
          system call fastpath
          lll unlock wake
          start_thread
      - 31.07% futex_wait
          do futex
          sys futex
          system_call_fastpath
          __lll_lock_wait
          start_thread
     + 0.01% futex_wait_setup
     + 0.01% futex_wake
                    [kernel.kallsyms]
                                        [k] raw spin unlock irgrestore
           a.out
                    [kernel.kallsyms]
                                        [k] finish_task_switch
     4.00% a.out
     1.99%
                    libpthread-2.19.so [.] pthread mutex lock
           a.out
                                        [.] pthread mutex unlock
     1.86% a.out
                    libpthread-2.19.so
```



# Perf Commands – perf annotate

perf record를 통해 생성된 file을 instruction level에서 분석한다.

\$ sudo perf annotate \_raw\_spin\_lock

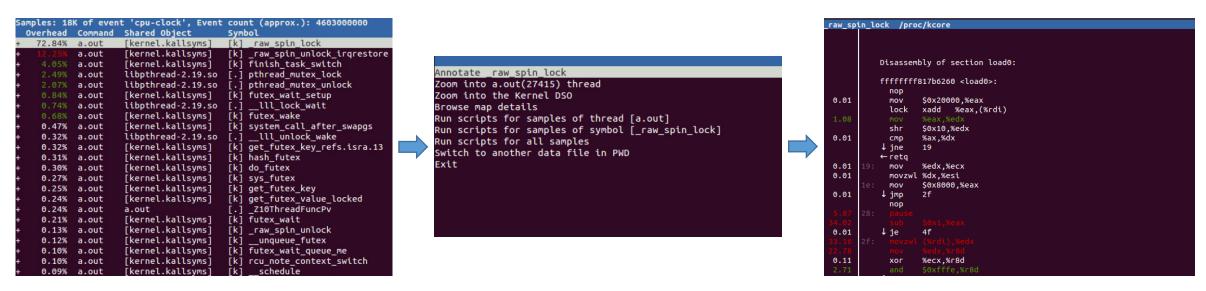
```
_raw_spin_lock /proc/kcore
           Disassembly of section load0:
           ffffffff817b6260 <load0>:
 0.01
                     $0x20000, %eax
                   xadd %eax,(%rdi)
                     $0x10,%edx
 0.01
                     %ax,%dx
            J ine
           ← retq
 0.01
                     %edx,%ecx
 0.01
             movzwl %dx,%esi
                     $0x8000,%eax
 0.01
            J jmp
            ↓ je
 0.01
                    %ecx,%r8d
 0.11
                     $0xfffe,%r8d
```



## Perf Commands – perf annotate

perf record를 통해 생성된 file을 instruction level에서 분석한다.

perf report 화면에서 annotate command 수행 가능



오른쪽 화살표키(->) 누름

엔터키 누름



### More about Perf...

그 밖의 다양한 방식의 측정이 가능하다.

Perf Tutorial https://perf.wiki.kernel.org/index.php/Tutorial



### Thank You

