# yogini

*Stretching the Linux Scheduler...*
*...to its Limits*

## Len Brown
### Sr. Principal Engineer
### Intel Open Source Technology Center

# Opportunity

A tool integrating...

1. workload generation
2. hardware and software observation
3. report generation

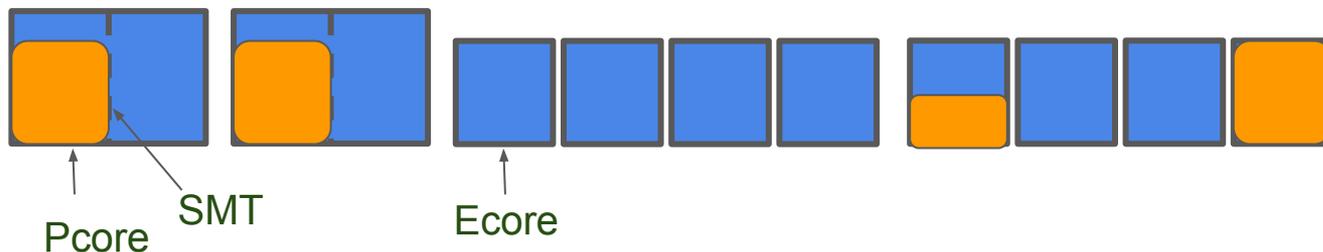Useful for scheduler+power+performance...

1. design
2. debug
3. tuning
4. regression testing

# Agenda

1.  Example
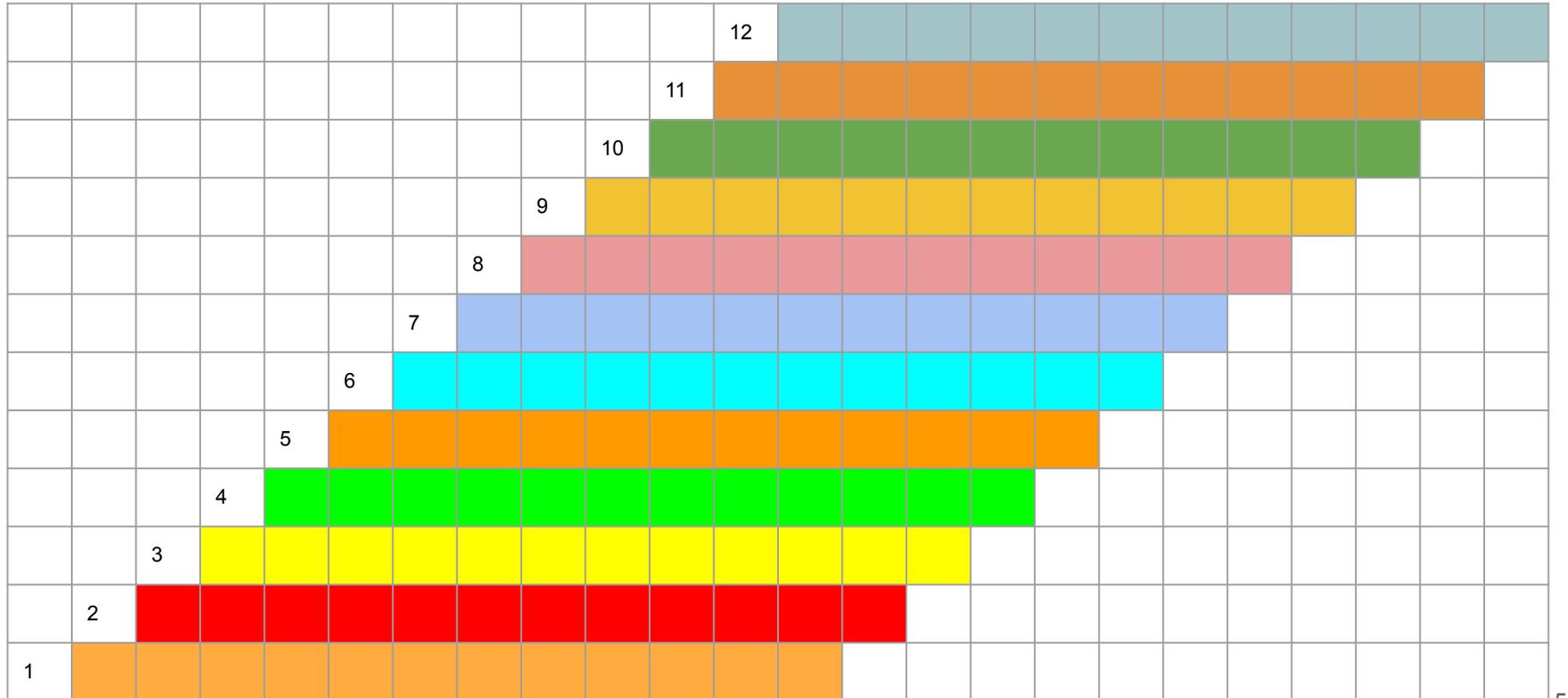2.  How yogini works
3.  Another Example

* Other names and brands may be claimed as property of others

# Linux v5.16 ITMT on Intel 2xPcore + 8xEcore

Task Placement:
1. Pcore
2. Ecore
3. Pcore HT sibling

Pcore
SMT
Ecore

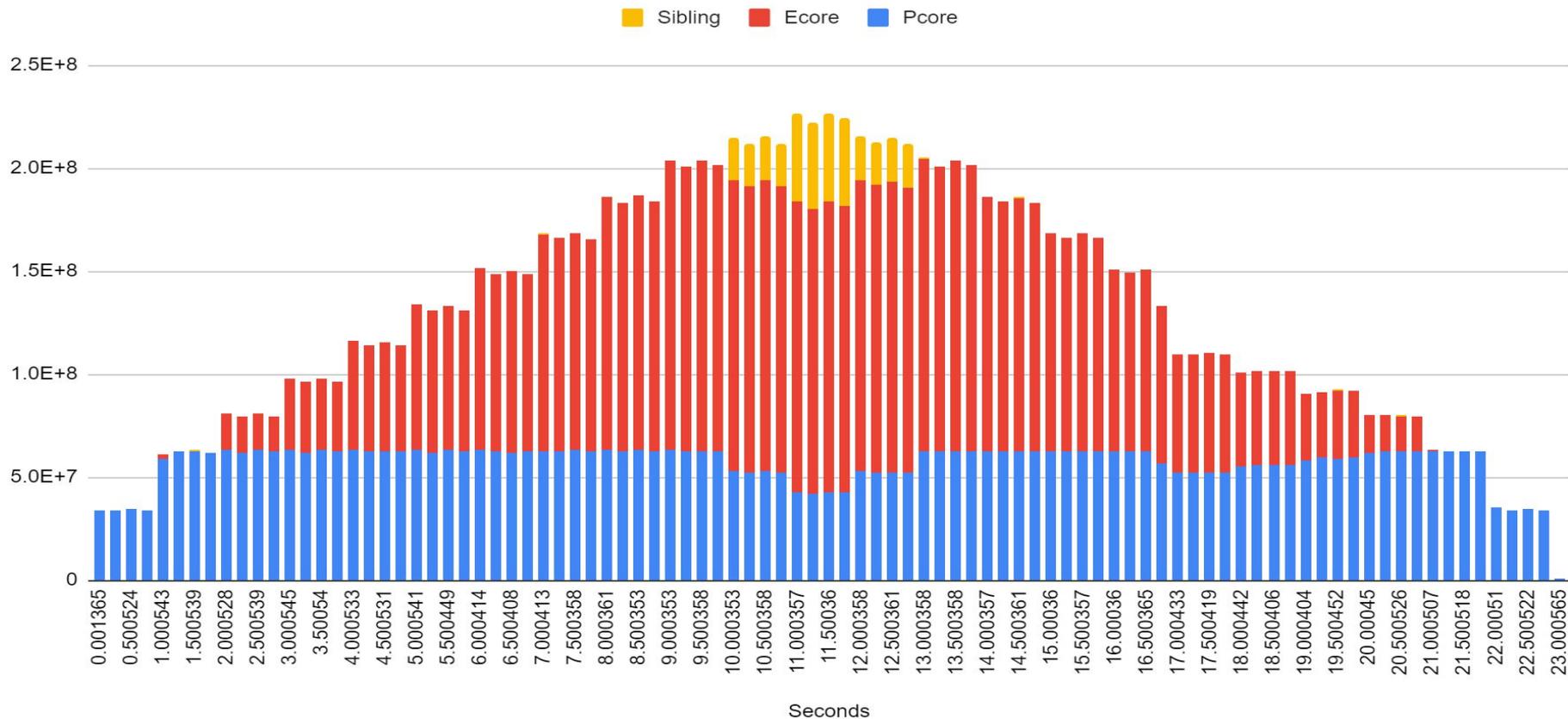Scheduler spreads to Ecore before HT sibling.
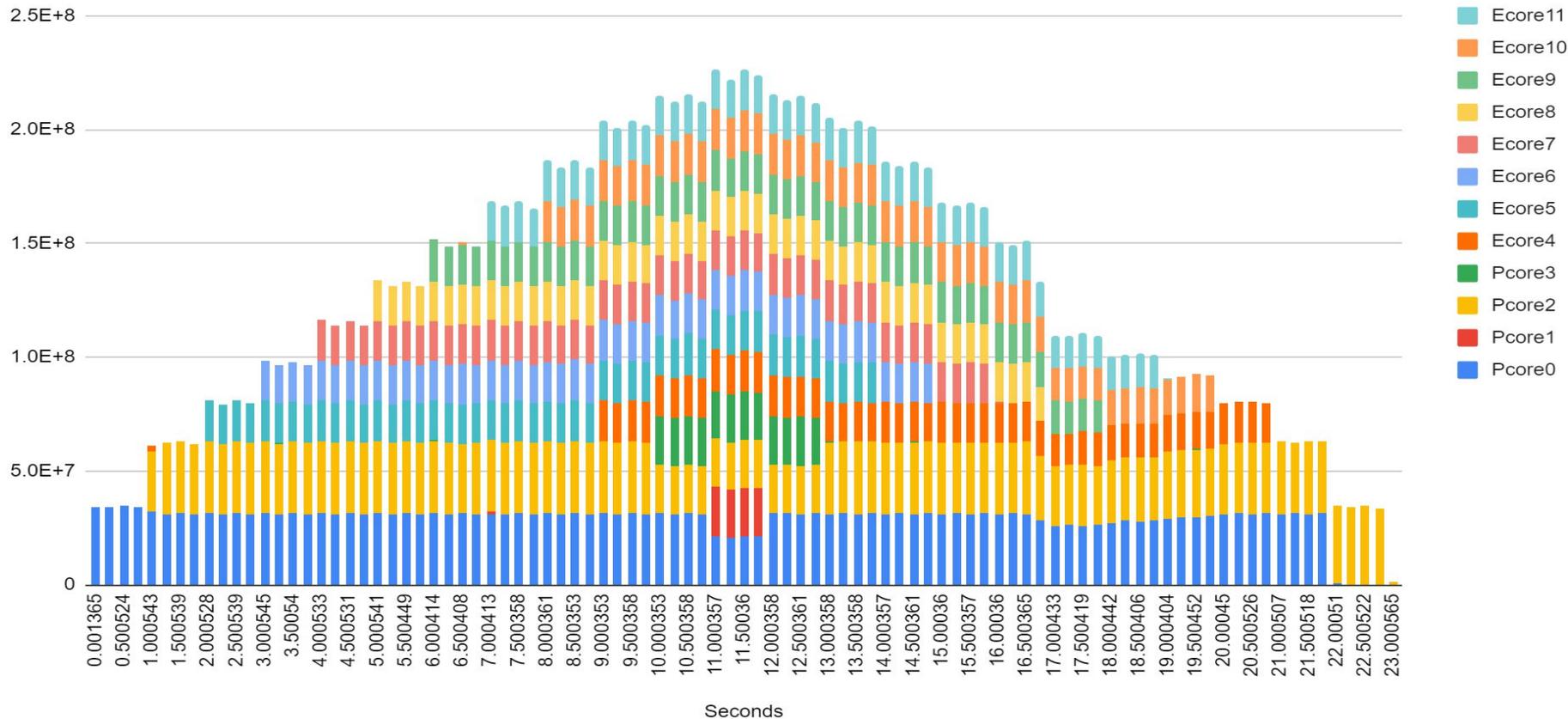
# 12-thread FIFO (100%) Stimulus

* Other names and brands may be claimed as property of others

# Yogini Pyramid100 Busy %  by CPU type



Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 work done by CPU type

Legend: Sibling (yellow), Ecore (red), Pcore (blue)

Y-axis: 2.5E+8, 2.0E+8, 1.5E+8, 5.0E+7, 0

X-axis (Seconds): 0.001365, 0.500524, 1.000543, 1.500539, 2.000528, 2.500539, 3.000545, 3.50054, 4.000533, 4.500531, 5.000541, 5.500449, 6.000414, 6.500408, 7.000413, 7.500358, 8.000361, 8.500353, 9.000353, 9.500358, 10.000353, 10.500358, 11.000357, 11.50036, 12.000358, 12.500361, 13.000358, 13.500358, 14.000357, 14.500361, 15.00036, 15.500357, 16.00036, 16.500365, 17.000433, 17.500419, 18.000442, 18.500406, 19.000404, 19.500452, 20.00045, 20.500526, 21.000507, 21.500518, 22.00051, 22.500522, 23.005565

Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Yogini pyramid100 work done per CPU over time



Linux Plumbers Conference    Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Yogini pyramid100 Frequency vs CPU over time



Legend:
- Pcore0
- Pcore1
- Pcore2
- Pcore3
- Ecore4
- Ecore5
- Ecore6
- Ecore7
- Ecore8
- Ecore9
- Ecore10
- Ecore11

X-axis: Seconds

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 Uncore MHz vs. time



Uncore Frequency [MHz]

Uncore Frequency [MHz]/Seconds

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 - IPC/CPU vs time



Legend: Pcore0, Pcore1, Pcore2, Pcore3, Ecore4, Ecore5, Ecore6, Ecore7, Ecore8, Ecore9, Ecore10, Ecore11

X-axis: Seconds

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 - IRQ/250ms per CPU



Legend: Pcore0, Pcore1, Pcore2, Pcore3, Ecore4, Ecore5, Ecore6, Ecore7, Ecore8, Ecore9, Ecore10, Ecore11

X-axis: Seconds

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 Temperature/CPU vs time



Legend: Pcore0, Pcore1, Pcore2, Pcore3, Ecore4, Ecore5, Ecore6, Ecore7, Ecore8, Ecore9, Ecore10, Ecore11

X-axis: Seconds

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini pyramid100 Volts/CPU vs time



Legend: Pcore0, Pcore1, Pcore2, Pcore3, Ecore4, Ecore5, Ecore6, Ecore7, Ecore8, Ecore9, Ecore10, Ecore11

X-axis: Seconds

* Other names and brands may be claimed as property of others

# Package, IA, UNCORE and GFX Power vs time



Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Yogini Pyramid100 Thread Work done by CPU



Legend:
- 12-GETCPU
- 11-GETCPU
- 10-GETCPU
- 9-GETCPU
- 8-GETCPU
- 7-GETCPU
- 6-GETCPU
- 5-GETCPU
- 4-GETCPU
- 3-GETCPU
- 2-GETCPU
- 1-GETCPU

Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Yogini Pyramid100 Thread CPU Residency Trace



Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Yogini Purpose:  The ability to easily...

1.  Generate well-understood workloads, to challenge Linux PM & scheduler
2.  Observe scheduler's success/failure against those challenges
3.  Foundation for regression test suite, to assure continuous improvement

\* Other names and brands may be claimed as property of others

# Yogini Goals: It needs to be easy to...

1. Install
2. Run on any topology
3. Run on any version of Linux
4. Share results
5. Understand results
6. Reproduce results
7. Compare before/after
8. Extend with additional workloads

* Other names and brands may be claimed as property of others

# Quick Start: Install, Run, Observe, Share

```
# tar zxf yogini-VERSION.tar.gz

# cd yogini-VERSION

# ./yogini > output.tsv
```

google sheets: Import output.tsv

         select data region, click "Insert Chart"

         click SHARE

# Optional Worker Parameters

Worker type

Number of copies (threads)

Waveform: eg. Rate of work, @ begin, @end

Start time, end-time, duty-cycle

affinity: start, stop, permanent

```
man ./yogini.8
```

# How yogini works

1. Calibrate Hardware

2. Start System Monitor

3. Run work

4. Output Results

Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Calibration sets "100%: Performance

For every workload type, in a test, 100% performance must be known

1. use pre-calibrate: `–calibrate AVX,12345678`

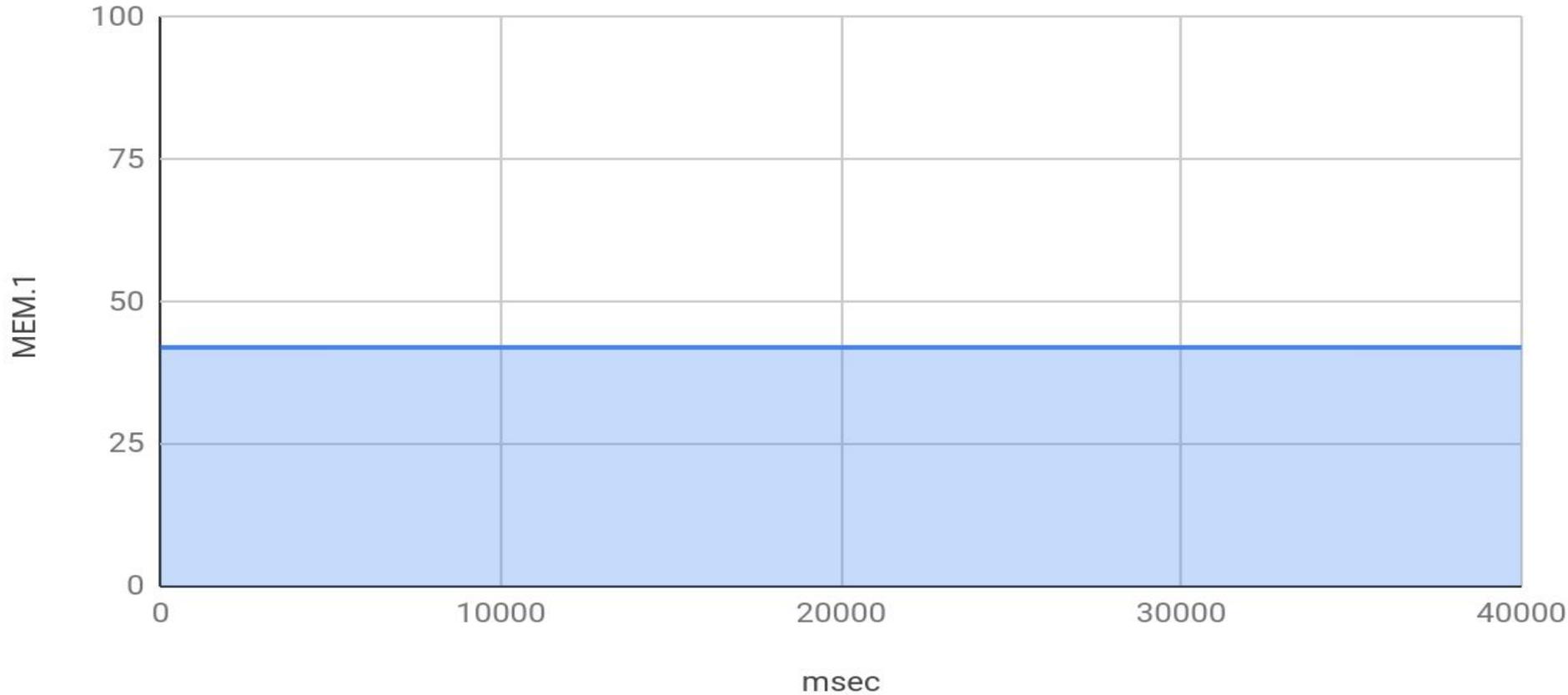2. measure on cpu0, or fastest of N CPUs: `–calibrate N`

# yogini -w rate100



Linux Plumbers Conference   Dublin, Ireland Sept 2022

# yogini -w start-msec1000,stop-msec3500



Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# yogini -w duty-cycle50

yogini -w rate42

* Other names and brands may be claimed as property of others

# yogini -w rate42,duty-cycle50



Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Constant Rate of Work/Time

# yogini -w rate1-100



Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Variable Work/Time

* Other names and brands may be claimed as property of others

# Work != Utilization

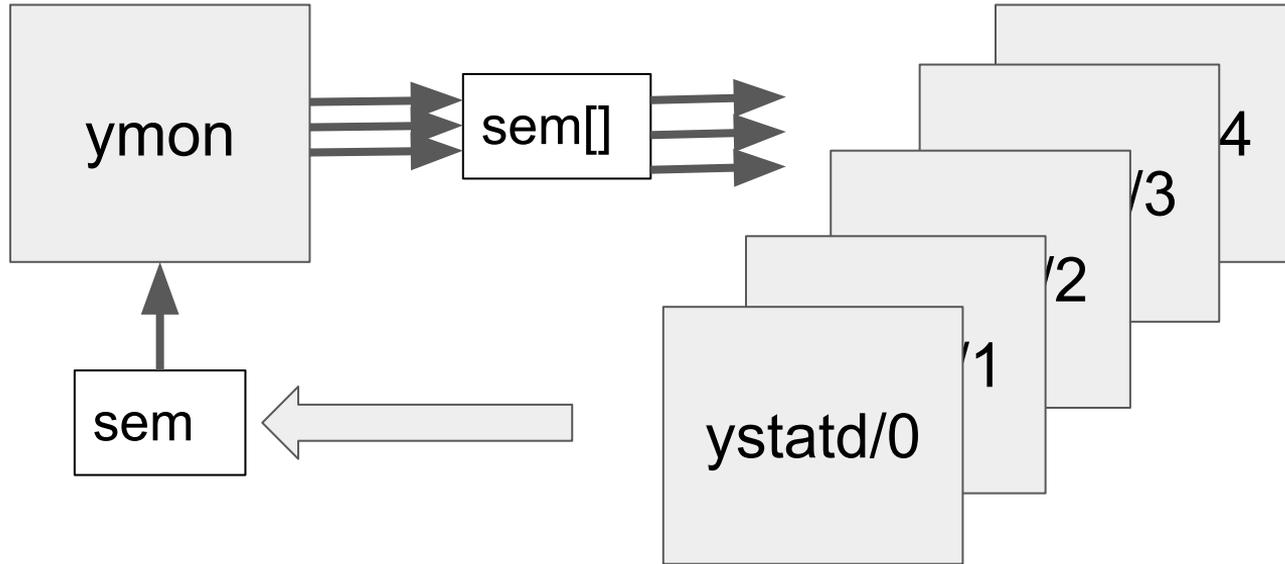Due to opportunistic turbo, 100% is often not attainable, or sustainable.

# Yogini system monitor

monitor thread periodically* collects:

1. Utilization Busy% (per CPU)
2. RunQ length (per CPU)
3. Frequency (per CPU)
4. Linux run queue length (per CPU)
5. IRQs (per CPU)
6. Instructions per Cycle (IPC)
7. Temperature (per CPU DTS)
8. RAPL power (package, CPU, GFX, RAM, Uncore)

\* `--monitor wakemsec250` (default 250 msec)

# system monitor architecture

# system monitor use

Monitor system for 10 sec:

```
# yogini
```

Fork my_program, monitor system until it exits:

```
# yogini my_program
```

Run built in AVX workload, monitor* until it exits

```
# yogini -w AVX
```

                                  * skip monitor with `--monitor off`

# Library of Built in Workloads

```
# yogini -w $WORKLOAD
```

WORKLOAD in:

1. GETCPU, RDTSC, PAUSE, TPAUSE, regAVX2, regVNNI

2. SSE, AVX, AVX2, AVX512, DOTPROD, VNNI

3. MEM, memcpy

# Working Set Size

GETCPU, RDTSC, regAVX2, regVNNI, PAUSE, TPAUSE **[No size]**

SSE, AVX, AVX2, AVX512, DOTPROD, VNNI **[ L1 dcache]**

MEM, memcpy **[L3 cache]**

Set working set size:

```
# yogini -w 256KB,AVX2 -w 100MB,MEM
```

# worker thread instrumentation

Worker thread time slide granularity [16.66 ms]

Self record every time slice:

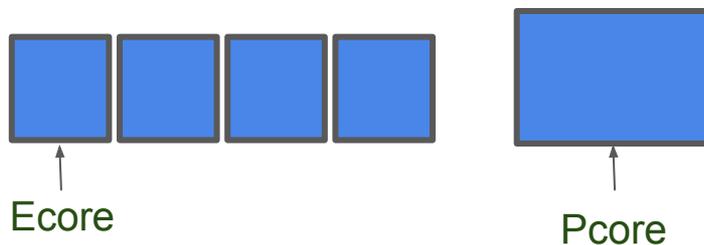1.   current CPU via getcpu(2)
2.   work-done counter

Set worker thread granularity to 10ms

```
# yogini -w wake-msec10
```

# Linux EAS test on 4xEcore + 1xPcore

Task Placement:
1. Pcore
2. Ecore
3. Pcore HT sibling

Ecore

Pcore

EAS: Ecores more efficient than Pcores at low MHz

# Example Ramp-Down on Big-Little

```
# yogini -w rate100-1
```
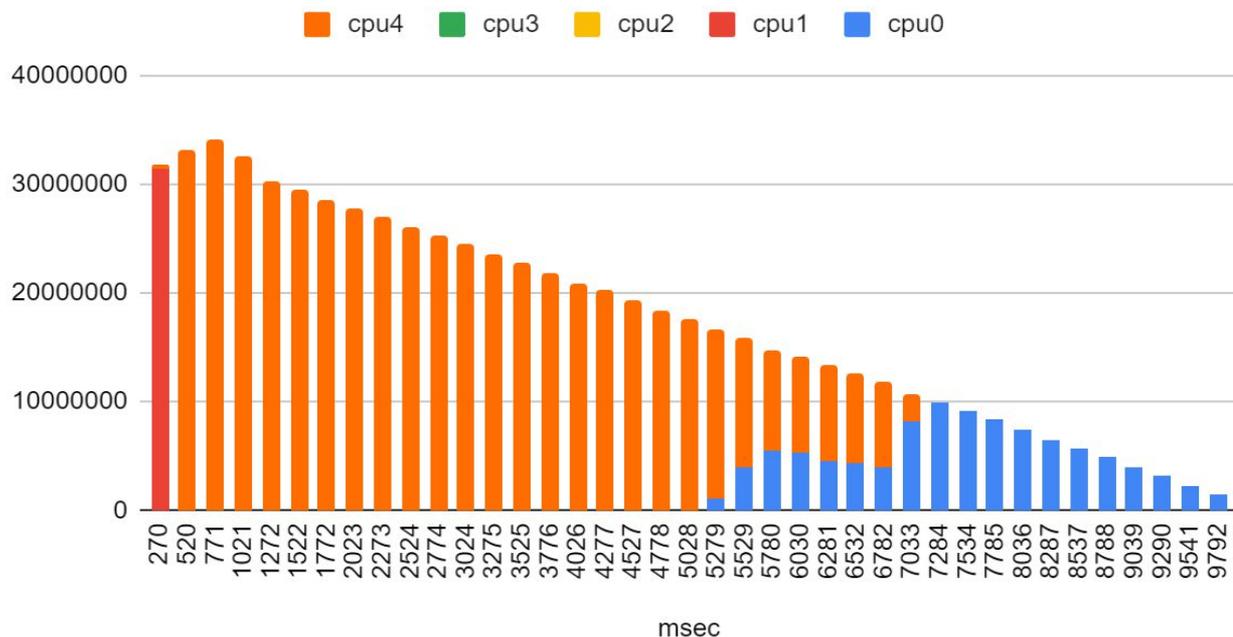
One thread

Requests 100% capacity, ramping down to 1%

Energy model marks Pcore4 as less efficient.

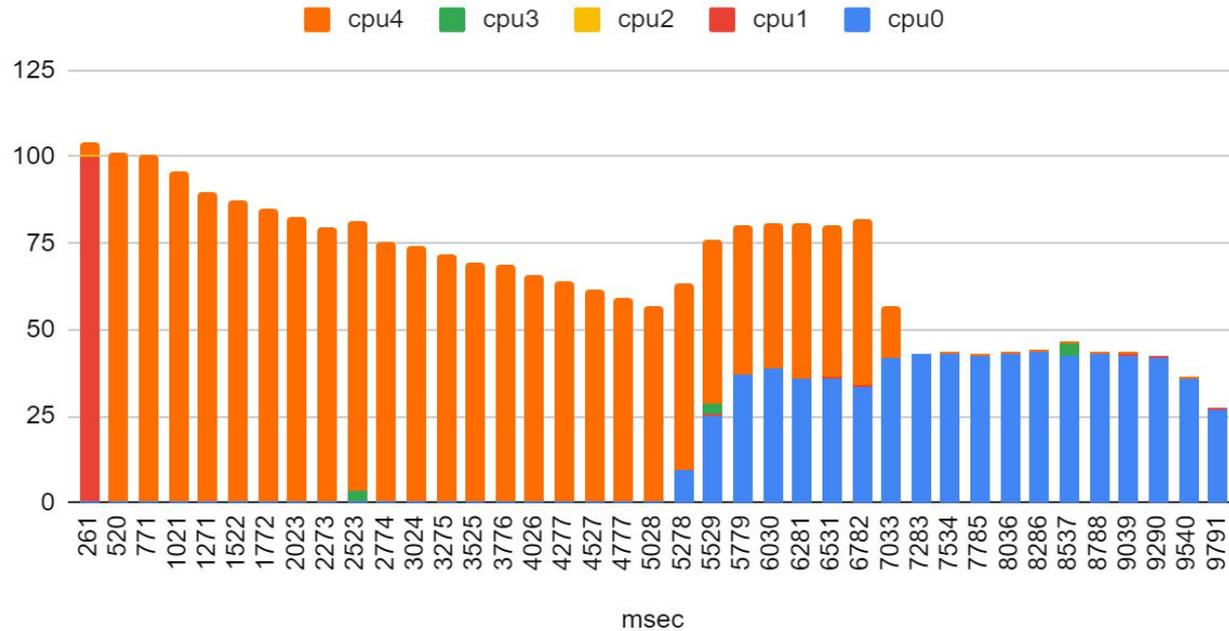Expect:  high demand to run on Pcore4, migrate to Ecore upon low demand

# Example Ramp-Down (Big -> Little) - Work Done



Work done per CPU vs Time

Legend: cpu4, cpu3, cpu2, cpu1, cpu0

Linux Plumbers Conference   Dublin, Ireland Sept 2022

# Example Ramp-Down (Big -> Little) %Busy



Busy% per CPU vs Time

Linux Plumbers Conference   Dublin, Ireland Sept 2022

* Other names and brands may be claimed as property of others

# Example Ramp Down - Frequency

Frequency per CPU vs Time

* Other names and brands may be claimed as property of others

# Example Ramp Down Temperature

## Temperature per CPU vs Time



Legend: cpu0, cpu1, cpu2, cpu3, cpu4

x-axis: msec

# Example Ramp Down - Power

RAPL Package Power vs Time

# Example Ramp Down: Summary

Summary Report:

- 100 Percent of Requested Throughput Achieved.
- 3.12 Average Watts
- 52 Task Migrations detected

Subjective Observations:

- Small->Big transition could have been faster
- Big -> small transition went meta-stable, but eventually worked

# What's Next?

What workloads are "interesting"?

Regression test scenarios?

Is .tsv the ultimate output?

Best way to distribute?